# Joint Universal Lossy Coding and Identification of Stationary Mixing Sources with General Alphabets

Maxim Raginsky, *Member, IEEE*

*Abstract*—We consider the problem of joint universal variable-rate lossy coding and identification for parametric classes of stationary $\beta$-mixing sources with general (Polish) alphabets. Compression performance is measured in terms of Lagrangians, while identification performance is measured by the variational distance between the true source and the estimated source. Provided that the sources are mixing at a sufficiently fast rate and satisfy certain smoothness and Vapnik–Chervonenkis learnability conditions, it is shown that, for bounded metric distortions, there exist universal schemes for joint lossy compression and identification whose Lagrangian redundancies converge to zero as $\sqrt{V_n \log n/n}$ as the block length $n$ tends to infinity, where $V_n$ is the Vapnik–Chervonenkis dimension of a certain class of decision regions defined by the $n$-dimensional marginal distributions of the sources; furthermore, for each $n$, the decoder can identify $n$-dimensional marginal of the active source up to a ball of radius $O(\sqrt{V_n \log n/n})$ in variational distance, eventually with probability one. The results are supplemented by several examples of parametric sources satisfying the regularity conditions.

*Index Terms*—Learning, minimum-distance density estimation, two-stage codes, universal vector quantization, Vapnik–Chervonenkis dimension.

## I. INTRODUCTION

It is well known that lossless source coding and statistical modeling are complementary objectives. This fact is captured by the Kraft inequality (see Section 5.2 in Cover and Thomas [1]), which provides a correspondence between uniquely decodable codes and probability distributions on a discrete alphabet. If one has full knowledge of the source statistics, then one can design an optimal lossless code for the source, and *vice versa*. However, in practice it is unreasonable to expect that the source statistics are known precisely, so one has to design *universal* schemes that perform asymptotically optimally within a given class of sources. In universal coding, too, as Rissanen has shown in [2], [3], the coding and modeling objectives can be accomplished jointly: given a sufficiently regular parametric family of discrete-alphabet sources, the encoder can acquire the source statistics via maximum-likelihood estimation on a sufficiently long data sequence and use this knowledge to select an appropriate coding scheme. Even in

nonparametric settings (e.g., the class of all stationary ergodic discrete-alphabet sources), universal schemes such as Ziv–Lempel [4] amount to constructing a probabilistic model for the source. In the reverse direction, Kieffer [5] and Merhav [6], among others, have addressed the problem of statistical modeling (parameter estimation or model identification) via universal lossless coding.

Once we consider *lossy* coding, though, the relationship between coding and modeling is no longer so simple. On the one hand, having full knowledge of the source statistics is certainly helpful for designing optimal rate-distortion codebooks. On the other hand, apart from some special cases (e.g., for i.i.d. Bernoulli sources and the Hamming distortion measure or for i.i.d. Gaussian sources and the squared-error distortion measure), it is not at all clear how to extract a reliable statistical model of the source from its reproduction via a rate-distortion code (although, as shown recently by Weissman and Ordentlich [7], the joint empirical distribution of the source realization and the corresponding codeword of a "good" rate-distortion code converges to the distribution solving the rate-distortion problem for the source). This is not a problem when the emphasis is on compression, but there are situations in which one would like to compress the source and identify its statistics at the same time. For instance, in *indirect adaptive control* (see, e.g., Chapter 7 of Tao [8]) the parameters of the plant (the controlled system) are estimated on the basis of observation, and the controller is modified accordingly. Consider the discrete-time stochastic setting, in which the plant state sequence is a random process whose statistics are governed by a finite set of parameters. Suppose that the controller is geographically separated from the plant and connected to it via a noiseless digital channel whose capacity is $R$ bits per use. Then, given the time horizon $T$, the objective is to design an encoder and a decoder for the controller to obtain reliable estimates of both the plant parameters and the plant state sequence from the $2^{TR}$ possible outputs of the decoder.

To state the problem in general terms, consider an information source emitting a sequence $\boldsymbol{X} = \{X_i\}_{i \in \mathbb{Z}}$ of random variables taking values in an alphabet $\mathcal{X}$. Suppose that the process distribution of $\boldsymbol{X}$ is not specified completely, but it is known to be a member of some parametric class $\{P_\theta : \theta \in \Lambda\}$. We wish to answer the following two questions:

1) Is the class $\{P_\theta : \theta \in \Lambda\}$ universally encodable with respect to a given single-letter distortion measure $\rho$, by codes with a given structure (e.g., all fixed-rate block

codes with a given per-letter rate, all variable-rate block codes, etc.)? In other words, does there exist a scheme that is asymptotically optimal for each $P_\theta$, $\theta \in \Lambda$?

2) If the answer to Question 1) is positive, can the codes be constructed in such a way that the decoder can not only reconstruct the source, but also identify its process distribution $P_\theta$, in an asymptotically optimal fashion?

In previous work [9], [10], we have addressed these two questions in the context of fixed-rate lossy block coding of stationary memoryless (i.i.d.) continuous-alphabet sources with parameter space $\Lambda$ a bounded subset of $\mathbb{R}^k$ for some finite $k$. We have shown that, under appropriate regularity conditions on the distortion measure and on the source models, there exist joint universal schemes for lossy coding and source identification whose redundancies (that is, the gap between the actual performance and the theoretical optimum given by the Shannon distortion-rate function) and source estimation fidelity both converge to zero as $O\big(\sqrt{\log n/n}\big)$, as the block length $n$ tends to infinity. The code operates by coding each block with the code matched to the source with the parameters estimated from the preceding block. Comparing this convergence rate to the $\log n/n$ convergence rate, which is optimal for redundancies of fixed-rate lossy block codes [11], we see that there is, in general, a price to be paid for doing compression and identification simultaneously. Furthermore, the constant hidden in the $O(\cdot)$ notation increases with the "richness" of the model class $\{P_\theta : \theta \in \Lambda\}$, as measured by the Vapnik–Chervonenkis (VC) dimension [12] of a certain class of measurable subsets of the source alphabet associated with the sources.

The main limitation of the results of [9], [10] is the i.i.d. assumption, which is rather restrictive as it excludes many practically relevant model classes (e.g., autoregressive sources, or Markov and hidden Markov processes). Furthermore, the assumption that the parameter space $\Lambda$ is bounded may not always hold, at least in the sense that we may not know the diameter of $\Lambda$ *a priori*. In this paper we relax both of these assumptions and study the existence and the performance of universal schemes for joint lossy coding and identification of stationary sources satisfying a mixing condition, when the sources are assumed to belong to a parametric model class $\{P_\theta : \theta \in \Lambda\}$, $\Lambda$ being an open subset of $\mathbb{R}^k$ for some finite $k$. Because the parameter space is not bounded, we have to use variable-rate codes with countably infinite codebooks, and the performance of the code is assessed by a composite Lagrangian functional [13] which captures the trade-off between the expected distortion and the expected rate of the code. Our result is that, under certain regularity conditions on the distortion measure and on the model class, there exist universal schemes for joint lossy source coding and identification such that, as the block length $n$ tends to infinity, the gap between the actual Lagrangian performance and the optimal Lagrangian performance achievable by variable-rate codes at that block length, as well as the source estimation fidelity at the decoder, converge to zero as $O\big(\sqrt{V_n \log n/n}\big)$, where $V_n$ is the VC dimension of a certain class of decision regions induced by the collection $\{P_\theta^n : \theta \in \Lambda\}$ of the $n$-

dimensional marginals of the source process distributions.

This result shows very clearly that the price to be paid for universality, in terms of both compression and identification, grows with the richness of the underlying model class, as captured by the VC dimension sequence $\{V_n\}$. The richer the model class, the harder it is to learn, which affects the compression performance of our scheme because we use the source parameters learned from past data to decide how to encode the current block. Furthermore, comparing the rate at which the Lagrangian redundancy decays to zero under our scheme with the $O(\log n/n)$ result of Chou, Effros and Gray [14], whose universal scheme is not aimed at identification, we immediately see that, in ensuring to satisfy the twin objectives of compression and modeling, we inevitably sacrifice some compression performance.

The paper is organized as follows. Section II introduces notation and basic concepts related to sources, codes and Vapnik–Chervonenkis classes. Section III lists and discusses the regularity conditions that have to be satisfied by the source model class, and contains the statement of our result. The result is proved in Section IV. Next, in Section V we give three examples of parametric source families (namely, i.i.d. Gaussian sources, Gaussian autoregressive sources and hidden Markov processes) which fit the framework of this paper under suitable regularity conditions. We conclude in Section VI and outline directions for future research. Finally, the Appendix contains some technical results on Lagrange-optimal variable-rate quantizers.

## II. Preliminaries

### A. Sources

In this paper, a *source* is a discrete-time stationary ergodic random process $\mathbf{X} = \{X_i\}_{i \in \mathbb{Z}}$ with alphabet $\mathcal{X}$. We assume that $\mathcal{X}$ is a Polish space (i.e., a complete separable metric space[1]) and equip $\mathcal{X}$ with its Borel $\sigma$-field. For any pair of indices $i, j \in \mathbb{Z}$ with $i < j$, let $X_i^j$ denote the segment $(X_i, X_{i+1}, \ldots, X_j)$ of $\mathbf{X}$. If $P$ is the process distribution of $\mathbf{X}$, then we let $\mathbb{E}_P\{\cdot\}$ denote expectation with respect to $P$, and let $P^n$ denote the marginal distribution of $X_1^n$. Whenever $P$ carries a subscript, e.g., $P = P_\theta$, we write $\mathbb{E}_\theta\{\cdot\}$ instead. We assume that there exists a fixed $\sigma$-finite measure $\mu$ on $\mathcal{X}$, such that the $n$-dimensional marginal of any process distribution of interest is absolutely continuous with respect to the product measure $\mu^n$, for all $n \geq 1$. We denote the corresponding densities $dP^n/d\mu^n$ by $p^n$. To avoid notational clutter, we omit the superscript $n$ from $\mu^n$, $P^n$ and $p^n$ whenever it is clear from the argument, as in $d\mu(x^n)$, $dP(x^n)$ or $p(x^n)$.

Given two probability measures $P, Q$ on a measurable space $(\mathcal{Z}, \mathcal{A})$, the *variational distance* between them is defined by

$$d(P, Q) \triangleq \sup_{\{A_i\} \subseteq \mathcal{A}} \sum_i |P(A_i) - Q(A_i)|,$$

where the supremum is over all finite $\mathcal{A}$-measurable partitions of $\mathcal{Z}$ (see, e.g., Section 5.2 of Gray [15]). If $p$ and $q$ are

---

[1]The canonical example is the Euclidean space $\mathbb{R}^d$ for some $d < \infty$.

the densities of $P$ and $Q$, respectively, with respect to a dominating measure $\nu$, then we can write

$$d(P,Q) = \int_{\mathcal{Z}} |p(z) - q(z)| d\nu(z).$$

A useful property of the variational distance is that, for any measurable function $f: \mathcal{Z} \to [0,1]$, $|\mathbb{E}_P f - \mathbb{E}_Q f| \leq d(P,Q)$. When $P$ and $Q$ are $n$-dimensional marginals of $P_\theta$ and $P_{\theta'}$, respectively, i.e., $P = P_\theta^n$ and $Q = P_{\theta'}^n$, we write $d_n(\theta, \theta')$ for $d(P_\theta^n, P_{\theta'}^n)$. If $\mathcal{A}'$ is a $\sigma$-subfield of $\mathcal{A}$, we define the variational distance $d(P,Q; \mathcal{A}')$ between $P$ and $Q$ with respect to $\mathcal{A}'$ by

$$d(P,Q;\mathcal{A}') \triangleq \sup_{\{A_i\} \subseteq \mathcal{A}'} \sum_i |P(A_i) - Q(A_i)|,$$

where the supremum is over all finite $\mathcal{A}'$-measurable partitions of $\mathcal{Z}$. Given a $\delta > 0$ and a probability measure $P$, the *variational ball* of radius $\delta$ around $P$ is the set of all probability measures $Q$ with $d(P,Q) \leq \delta$.

Given a source $\boldsymbol{X}$ with process distribution $P$, let $P_{-\infty}^0$ and $P_1^\infty$ denote the marginal distributions of $P$ on $\{X_i\}_{i \leq 0}$ and $\{X_i\}_{i \geq 1}$, respectively. For each $k \geq 1$, the *kth-order absolute regularity coefficient* (or $\beta$-*mixing coefficient*) of $P$ is defined as [16], [17]:

$$\beta_P(k) \triangleq \sup \left\{ \sum_i \sum_j |P(A_i \cap B_j) - P_{-\infty}^0(A_i)P_1^\infty(B_j)| \right\},$$

where the supremum is over all finite $\sigma(X_{-\infty}^0)$-measurable partitions $\{A_i\}$ and all finite $\sigma(X_k^\infty)$-measurable partitions $\{B_j\}$. Observe that

$$\beta_P(k) = d\left(P, P_{-\infty}^0 \times P_1^\infty; \sigma(X_{-\infty}^0, X_k^\infty)\right), \qquad (1)$$

the variational distance between $P$ and the product distribution $P_{-\infty}^0 \times P_1^\infty$ with respect to the $\sigma$-algebra $\sigma(X_{-\infty}^0, X_k^\infty)$. Since $\boldsymbol{X}$ is stationary, we can "split" its process distribution at any point $l \in \mathbb{Z}$ and define $\beta_P(k)$ equivalently by

$$\beta_P(k) \triangleq d\left(P, P_{-\infty}^l \times P_{l+1}^\infty; \sigma(X_{-\infty}^l, X_{l+k}^\infty)\right). \qquad (2)$$

Again, if $P$ is subscripted by some $\theta$, $P = P_\theta$, then we write $\beta_\theta(k)$.

*B. Codes*

The class of codes we consider here is the collection of all finite-memory variable-rate vector quantizers. Let $\widehat{\mathcal{X}}$ be a *reproduction alphabet*, also assumed to be Polish. We assume that $\mathcal{X} \cup \widehat{\mathcal{X}}$ is a subset of a Polish metric space $\mathcal{Y}$ with a bounded metric $\rho_0(\cdot,\cdot)$: there exists some $\rho_{\max} < +\infty$, such that $\rho_0(y, y') \leq \rho_{\max}$ for all $y, y' \in \mathcal{Y}$. We take $\rho: \mathcal{X} \times \widehat{\mathcal{X}} \to [0, \rho_{\max}]$, $\rho(x, \widehat{x}) \triangleq \rho_0(x, \widehat{x})$, as our (single-letter) *distortion function*. A variable-rate vector quantizer with block length $n$ and memory length $m$ is a pair $C^{n,m} = (f, \varphi)$, where $f: \mathcal{X}^n \times \mathcal{X}^m \to \mathcal{S}$ is the *encoder*, $\varphi: \mathcal{S} \to \widehat{\mathcal{X}}^n$ is the *decoder*, and $\mathcal{S} \subseteq \{0,1\}^*$ is a countable collection of binary strings satisfying the prefix condition or, equivalently, the Kraft inequality

$$\sum_{s \in \mathcal{S}} 2^{-\ell(s)} \leq 1,$$

where $\ell(s)$ denotes the length of $s$ in bits. The mapping of the source $\boldsymbol{X}$ into the reproduction process $\widehat{\boldsymbol{X}}$ is defined by

$$\widehat{X}_{nk+1}^{n(k+1)} = \varphi\left(f\left(X_{nk+1}^{n(k+1)}, X_{nk-m+1}^{nk}\right)\right), \qquad k \in \mathbb{Z}.$$

That is, the encoding is done in blocks of length $n$, but the encoder is also allowed to observe the $m$ symbols immediately preceding each block. The *effective memory* of $C^{n,m}$ is defined as the set $\mathcal{M} \subseteq \{1, \ldots, m\}$, such that

$$f(x^m) = f(\widetilde{x}^m), \qquad \forall x^m, \widetilde{x}^m \in \mathcal{X}^m : x_i = \widetilde{x}_i, \forall i \in \mathcal{M}.$$

The size $|\mathcal{M}|$ of $\mathcal{M}$ is called the *effective memory length* of $C^{n,m}$. We shall often use $C^{n,m}$ to also denote the composite mapping $\varphi \circ f$: $\widehat{X}_1^n = C^{n,m}(X_1^n, X_{-m+1}^0)$. When the code has zero memory ($m = 0$), we shall denote it more compactly by $C^n$.

The performance of the code on the source with process distribution $P$ is measured by its expected distortion

$$D_P(C^{n,m}) \triangleq \mathbb{E}_P\left\{\rho_n(X_1^n, \widehat{X}_1^n)\right\},$$

where for $x^n \in \mathcal{X}^n$ and $\widehat{x}^n \in \widehat{\mathcal{X}}^n$, $\rho_n(x^n, \widehat{x}^n) \triangleq n^{-1} \sum_{i=1}^n \rho(x_i, \widehat{x}_i)$ is the per-letter distortion incurred in reproducing $x^n$ by $\widehat{x}^n$, and by its expected rate

$$R_P(C^{n,m}) \triangleq \mathbb{E}_P\left\{\ell_n\left(f\left(X_1^n, X_{-m+1}^0\right)\right)\right\},$$

where $\ell_n(s)$ denotes the length of a binary string $s$ in bits, normalized by $n$. (We follow Neuhoff and Gilbert [18] and normalize the distortion and the rate by the length $n$ of the *reproduction* block, not by the combined length $n + m$ of the source block plus the memory input.) When working with variable-rate quantizers, it is convenient [13], [19] to absorb the distortion and the rate into a single performance measure, the *Lagrangian distortion*

$$L_P(C^{n,m}, \lambda) \triangleq D_P(C^{n,m}) + \lambda R_P(C^{n,m}),$$

where $\lambda > 0$ is the *Lagrange multiplier* which controls the distortion-rate trade-off. Geometrically, $L_P(C^{n,m})$ is the $y$-intercept of the line with slope $-\lambda$, passing through the point $(R_P(C^{n,m}), D_P(C^{n,m}))$ in the rate-distortion plane [20]. If $P$ carries a subscript, $P = P_\theta$, then we write $D_\theta(\cdot)$, $R_\theta(\cdot)$ and $L_\theta(\cdot)$.

*C. Vapnik–Chervonenkis classes*

In this paper, we make heavy use of Vapnik–Chervonenkis theory (see Devroye, Györfi and Lugosi [21], Vapnik [22], Devroye and Lugosi [23] or Vidyasagar [24] for detailed treatments). This section contains a brief summary of the needed concepts and results. Let $(\mathcal{Z}, \mathcal{A})$ be a measurable space. For any collection $\mathcal{C} \subseteq \mathcal{A}$ of measurable subsets of $\mathcal{Z}$ and any $n$-tuple $z^n \in \mathcal{Z}^n$, define the set $\mathcal{C}(z^n) \subseteq \{0,1\}^n$ consisting of all distinct binary strings of the form $(1_{\{z_1 \in A\}}, \ldots, 1_{\{z_n \in A\}})$, $A \in \mathcal{C}$. Then

$$\mathsf{S}_n(\mathcal{C}) \triangleq \max_{z^n \in \mathcal{Z}^n} |\mathcal{C}(z^n)|$$

is called the *nth shatter coefficient* of $\mathcal{C}$. The *Vapnik–Chervonenkis dimension* (or VC-dimension) of $\mathcal{C}$, denoted by $\mathsf{V}(\mathcal{C})$, is defined as the largest $n$ for which $\mathsf{S}_n(\mathcal{C}) = 2^n$ (if

$S_n(\mathcal{C}) = 2^n$ for all $n = 1, 2, \ldots$, then we set $V(\mathcal{C}) = \infty$). If $V(\mathcal{C}) < \infty$, then $\mathcal{C}$ is called a *Vapnik–Chervonenkis class* (or VC class). If $\mathcal{C}$ is a VC class with $V(\mathcal{C}) \geq 2$, then it follows from the results of Vapnik and Chervonenkis [12] and Sauer [25] that $S_n(\mathcal{C}) \leq n^{V(\mathcal{C})}$.

For a VC class $\mathcal{C}$, the so-called *Vapnik–Chervonenkis inequalities* (see Lemma 2.1 below) relate its VC dimension $V(\mathcal{C})$ to maximal deviations of the probabilities of the events in $\mathcal{C}$ from their relative frequencies with respect to an i.i.d. sample of size $n$. For any $z^n \in \mathcal{Z}^n$, let

$$P_{z^n} = \frac{1}{n} \sum_{i=1}^{n} \delta_{z_i}$$

denote the induced empirical distribution, where $\delta_{z_i}$ is the Dirac measure (point mass) concentrated at $z_i$. We then have the following:

*Lemma 2.1 (Vapnik–Chervonenkis inequalities):* Let $P$ be a probability measure on $(\mathcal{Z}, \mathcal{A})$, and $Z_1^n = (Z_1, \ldots, Z_n)$ an $n$-tuple of independent random variables with $Z_i \sim P$, $1 \leq i \leq n$. Let $\mathcal{C}$ be a Vapnik–Chervonenkis class with $V(\mathcal{C}) \geq 2$. Then for every $\delta > 0$,

$$\mathbb{P}\left\{ \sup_{A \in \mathcal{C}} |P_{Z_1^n}(A) - P(A)| > \delta \right\} \leq 8 n^{V(\mathcal{C})} e^{-n\delta^2/32} \quad (3)$$

and

$$\mathbb{E}\left\{ \sup_{A \in \mathcal{C}} |P_{Z_1^n}(A) - P(A)| \right\} \leq c\sqrt{\frac{V(\mathcal{C}) \log n}{n}}, \quad (4)$$

where $c > 0$ is a universal constant. The probabilities and expectations are with respect to the product measure $P^n$ on $(\mathcal{Z}^n, \mathcal{A}^n)$.

*Remark 2.1:* A more refined technique involving metric entropies and empirical covering numbers, due to Dudley [26], can yield a much better $O(1/\sqrt{n})$ bound on the expected maximal deviation between the true and the empirical probabilities. This improvement, however, comes at the expense of a much larger constant hidden in the $O(\cdot)$ notation.

Finally, we shall need the following lemma, which is a simple corollary of the results of Karpinski and Macintyre [27] (see also Section 10.3.5 of Vidyasagar [24]):

*Lemma 2.2:* Let $\mathcal{C} = \{A_\xi : \xi \in \mathbb{R}^k\}$ be a collection of measurable subsets of $\mathbb{R}^d$, such that

$$A_\xi = \{z \in \mathbb{R}^d : \Pi(z, \xi) > 0\}, \qquad \xi \in \mathbb{R}^k$$

where for each $z \in \mathbb{R}^d$, $\Pi(z, \cdot)$ is a polynomial of degree $s$ in the components of $\xi$. Then $\mathcal{C}$ is a VC class with $V(\mathcal{C}) \leq 2k \log(4es)$.

## III. STATEMENT OF RESULTS

In this section we state our result concerning universal schemes for joint lossy compression and identification of stationary sources under certain regularity conditions. We work in the usual setting of universal source coding: we are given a source $\mathbf{X} = \{X_i\}_{i \in \mathbb{Z}}$ whose process distribution is known to be a member of some parametric class $\{P_\theta : \theta \in \Lambda\}$. The parameter space $\Lambda$ is an open subset of the Euclidean space $\mathbb{R}^k$ for some finite $k$, and we assume that $\Lambda$ has nonempty interior. We wish to design a sequence of variable-rate vector quantizers, such that the decoder can reliably reconstruct the original source sequence $\mathbf{X}$ and reliably identify the active source in an asymptotically optimal manner for all $\theta \in \Lambda$. We begin by listing the regularity conditions.

*Condition 1.* The sources in $\{P_\theta : \theta \in \Lambda\}$ are *algebraically $\beta$-mixing*: there exists a constant $r > 0$, such that

$$\beta_\theta(k) = O(k^{-r}), \qquad \forall \theta \in \Lambda$$

where the constant implicit in the $O(\cdot)$ notation may depend on $\theta$.

This condition ensures that certain finite-block functions of the source $\mathbf{X}$ can be approximated in distribution by i.i.d. processes, so that we can invoke the Vapnik–Chervonenkis machinery of Section II-C. This "blocking" technique, which we exploit in the proof of our Theorem 3.1, dates back to Bernstein [28], and was used by Yu [29] to derive rates of convergence in the uniform laws of large numbers for stationary mixing processes, and by Meir [30] in the context of nonparametric adaptive prediction of stationary time series. As an example of when an even stronger decay condition holds, let $\mathbf{X} = \{X_i\}_{i \in \mathbb{Z}}$ be a finite-order autoregressive moving-average (ARMA) process driven by a zero-mean i.i.d. process $\mathbf{Y} = \{Y_i\}$, i.e., there exist poisitive integers $p, q$ and $p + q + 1$ real constants $a_0, a_1 \ldots, a_p, b_1, \ldots, b_q$ such that

$$\sum_{i=0}^{p} a_i X_{n-i} = \sum_{j=1}^{q} b_j Y_{n-j}, \qquad n \in \mathbb{Z}.$$

Mokkadem [31] has shown that, provided the common distribution of the $Y_i$ is absolutely continuous and the roots of the polynomial $A(z) = \sum_{i=0}^{p} a_i z^i$ lie outside the unit circle in the complex plane, the $\beta$-mixing coefficients of $\mathbf{X}$ decay to zero exponentially.

*Condition 2.* For each $\theta \in \Lambda$, there exist constants $\delta_\theta > 0$ and $c_\theta > 0$, such that

$$\sup_n \frac{d_n(\theta, \theta')}{\sqrt{n}} \leq c_\theta \|\theta - \theta'\|$$

for all $\theta'$ in the open ball of radius $\delta_\theta$ centered at $\theta$, where $\|\cdot\|$ is the Euclidean norm on $\Lambda$.

This condition guarantees that, for any sequence $\{\delta_n\}_{n \in \mathbb{N}}$ of positive reals such that

$$\delta_n \to 0, \sqrt{n}\delta_n \to 0, \qquad \text{as } n \to \infty$$

and any sequence $\{\theta_n\}_{n \in \mathbb{N}}$ in $\Lambda$ satisfying $\|\theta_n - \theta\| < \delta_n$ for a given $\theta \in \Lambda$, we have

$$d_n(\theta, \theta_n) \to 0, \qquad \text{as } n \to \infty.$$

It is weaker (i.e., more general) than the conditions of Rissanen [2], [3] which control the behavior of the relative entropy (information divergence) as a function of the source parameters in terms of the Fisher information and related quantities. Indeed,

for each $n$ let

$$
\begin{aligned}
D_n(P_\theta \| P_{\theta'}) &= \frac{1}{n} \mathbb{E}_\theta \left\{ \ln \frac{dP_\theta}{dP_{\theta'}}(X_1^n) \right\} \\
&\equiv \frac{1}{n} \int_{\mathcal{X}^n} p_\theta(x^n) \ln \frac{p_\theta(x^n)}{p_{\theta'}(x^n)} d\mu(x^n)
\end{aligned}
$$

be the normalized $n$th-order relative entropy (information divergence) between $P_\theta$ and $P_{\theta'}$. Suppose that, for each $\theta$, $D_n(P_\theta \| P_{\theta'})$ is twice continuously differentiable as a function of $\theta'$. Let $\theta'$ lie in an open ball of radius $\delta$ around $\theta$. Since $D(P_\theta \| P_{\theta'})$ attains its minimum at $\theta' = \theta$, the gradient $\nabla_{\theta'} D_n(P_\theta \| P_{\theta'})$ evaluated at $\theta' = \theta$ is zero, and we can write the second-order Taylor expansion of $D_n$ about $\theta$ as

$$
D_n(P_\theta \| P_{\theta'}) = \frac{1}{2}(\theta - \theta')^T J_n(\theta)(\theta - \theta') + o(\|\theta - \theta'\|^2), \quad (5)
$$

where the Hessian matrix

$$
J_n(\theta) = \left[ \frac{\partial^2}{\partial \theta'_i \partial \theta'_j} D_n(P_\theta \| P_{\theta'}) \Big|_{\theta' = \theta} \right]_{i,j=1,\ldots,k},
$$

under additional regularity conditions, is equal to the Fisher information matrix

$$
I_n(\theta) = \left[ -\frac{1}{n} \mathbb{E}_\theta \left\{ \frac{\partial^2}{\partial \theta_i \partial \theta_j} \ln p_\theta(X_1^n) \right\} \right]_{i,j=1,\ldots,k}
$$

(see Clarke and Barron [32]). Assume now, following Rissanen [2], [3], that the sequence of matrix norms $\{\|I_n(\theta)\|\}$ is bounded (by a constant depending on $\theta$). Then we can write

$$
\begin{aligned}
D_n(P_\theta \| P_{\theta'}) &\leq \frac{1}{2}(\|I_n(\theta)\| + o(1)) \cdot \|\theta - \theta'\|^2 \\
&\leq c'_\theta \|\theta - \theta'\|^2,
\end{aligned}
$$

i.e., the normalized relative entropies $D_n(P_\theta \| P_{\theta'})$ are locally quadratic in $\theta'$. Then Pinsker's inequality (see, e.g., Lemma 5.2.8 of Gray [15]) implies that $\sqrt{2D_n(P_\theta \| P_{\theta'})} \geq d_n(\theta, \theta')/\sqrt{n}$, and we recover our Condition 2. Rissanen's condition, while stronger than our Condition 2, is easier to check, the fact which we exploit in our discussion of examples of Section V.

*Condition 3.* For each $n$, let $\mathcal{A}_n$ be the collection of all sets of the form

$$
A_{\theta,\theta'} = \left\{ x^n \in \mathcal{X}^n : p_\theta(x^n) > p_{\theta'}(x^n) \right\}, \qquad \theta \neq \theta'.
$$

Then we require that, for each $n$, $\mathcal{A}_n$ is a VC class, and $\mathsf{V}(\mathcal{A}_n) = o(n/\log n)$.

This condition is satisfied, for example, when $\mathsf{V}(\mathcal{A}_n) = V < \infty$ independently of $n$, or when $\mathsf{V}(\mathcal{A}_n) = \log n$. The use of the class $\mathcal{A}_n$ dates back to the work of Yatracos [33] on minimum-distance density estimation. The ideas of Yatracos were further developed by Devroye and Lugosi [34], [35], who dubbed $\mathcal{A}_n$ the *Yatracos class* (associated with the densities $p_\theta^n$). We shall adhere to this terminology. To give an intuitive interpretation to $\mathcal{A}_n$, let us consider a pair $\theta, \theta' \in \Lambda$ of distinct parameter vectors and note that the set $\{x^n : p_\theta(x^n) > p_{\theta'}(x^n)\}$ consists of all $x^n$ for which the simple hypothesis test

$$
H_0 : X_1^n \sim P_\theta^n \quad \text{versus} \quad H_1 : X_1^n \sim P_{\theta'}^n \qquad (6)
$$

is passed by the null hypothesis $H_0$ under the likelihood-ratio decision rule. Now, suppose that $Z_1, \ldots, Z_m$ are drawn *independently* from $P_\theta^n$. To each $A \in \mathcal{A}_n$ we can associate a *classifier* $\kappa_A : \mathcal{X}^n \to \{0, 1\}$ defined by $\kappa_A(x^n) \triangleq 1_{\{x^n \in A\}}$. Call two sets $A, A' \in \mathcal{A}_n$ *equivalent* with respect to the sample $Z_1^n = (Z_1, \ldots, Z_m)$, and write $A \sim_{Z_1^n} A'$, if their associated classifiers yield identical classification patterns:

$$
\big(\kappa_A(Z_1), \ldots, \kappa_A(Z_m)\big) = \big(\kappa_{A'}(Z_1), \ldots, \kappa_{A'}(Z_m)\big).
$$

It is easy to see that $\sim_{Z_1^n}$ is an equivalence relation. From the definitions of the shatter coefficients $\mathsf{S}_m(\mathcal{A}_n)$ and the VC dimension $\mathsf{V}(\mathcal{A}_n)$ (cf. Section II-C), we see that the cardinality of the quotient set $\mathcal{A}_n / \sim_{Z_1^n}$ is equal to $2^m$ for all sample sizes $m \leq \mathsf{V}(\mathcal{A}_n)$, whereas for $m > \mathsf{V}(\mathcal{A}_n)$, it is bounded from above by $m^{\mathsf{V}(\mathcal{A}_n)}$, which is strictly less than $2^m$. Thus, the fact that the Yatracos class $\mathcal{A}_n$ has finite VC dimension implies that the problem of estimating the density $p_\theta^n$ from a large i.i.d. sample reduces, in a sense, to a finite number (in fact, polynomial in the sample size $m$, for $m > \mathsf{V}(\mathcal{A}_n)$) of simple hypothesis tests of the type (6). Our Condition 1 will then allow us to transfer this intuition to (weakly) dependent samples.

Now that we have listed the regularity conditions that must hold for the sources in $\{P_\theta : \theta \in \Lambda\}$, we can state our main result.

*Theorem 3.1:* Let $\{P_\theta : \theta \in \Lambda\}$ be a parametric class of sources satisfying Conditions 1–3. Then for every $\lambda > 0$ and every $\eta > 0$, there exists a sequence $\{C_*^{n,m_n}\}_{n \in \mathbb{N}}$ of variable-rate vector quantizers with memory length $m_n \leq n(n + n^{(2+\eta)/r} + 1)$ and effective memory length $n^2$, such that, for all $\theta \in \Lambda$,

$$
L_\theta(C_*^{n,m_n}, \lambda) - \inf_{m \geq 0} \inf_{C^{n,m}} L_\theta(C^{n,m}, \lambda)
$$
$$
= O\left( \sqrt{\frac{\mathsf{V}(\mathcal{A}_n) \log n}{n}} \right), \quad (7)
$$

where the constants implicit in the $O(\cdot)$ notation depend on $\theta$. Furthermore, for each $n$, the binary description produced by the encoder is such that the decoder can identify the $n$-dimensional marginal of the active source up to a variational ball of radius $O\big(\sqrt{\mathsf{V}(\mathcal{A}_n) \log n / n}\big)$ with probability one.

What (7) says is that, for each block length $n$ and each $\theta \in \Lambda$, the code $C_*^{n,m_n}$, which is *independent of* $\theta$, performs almost as well as the best finite-memory quantizer with block length $n$ that can be designed with full *a priori* knowledge of the $n$-dimensional marginal $P_\theta^n$. Thus, as far as compression goes, our scheme can compete with all finite-memory variable-rate lossy block codes (vector quantizers), with the additional bonus of allowing the decoder to identify the active source in an asymptotically optimal manner.

It is not hard to see that the double infimum in (7) is achieved already in the zero-memory case, $m = 0$. Indeed, it is immediate that having nonzero memory can only improve the Lagrangian performance, i.e.,

$$
\inf_{m \geq 0} \inf_{C^{n,m}} L_\theta(C^{n,m}, \lambda) \leq \inf_{C^n} L_\theta(C^n, \lambda), \qquad \forall \theta \in \Lambda.
$$

On the other hand, given any code $C^{n,m} = (f, \varphi)$, we can construct a zero-memory code $C_0^n = (f_0, \varphi_0)$, such that $L_\theta(C_0^n, \lambda) \leq L_\theta(C^{n,m}, \lambda)$ for all $\theta \in \Lambda$. To see this, define for each $x^n \in \mathcal{X}^n$ the set

$$\mathcal{S}(x^n) \triangleq \{s \in \{0,1\}^* : s = f(x^n, z^m) \text{ for some } z^m \in \mathcal{X}^m\},$$

and let

$$f_0(x^n) = \underset{s \in \mathcal{S}(x^n)}{\arg\min} \left( \rho_n(x^n, \varphi(s)) + \lambda \ell(s) \right), \qquad \forall x^n \in \mathcal{X}^n$$

and $\varphi_0 \equiv \varphi$. Then, given any $(x^n, z^m) \in \mathcal{X}^n \times \mathcal{X}^m$, let $s = f(x^n, z^m)$. We then have

$$\begin{aligned}
\rho_n(x^n&, C_0^n(x^n)) + \lambda \ell(f_0(x^n)) \\
&= \rho_n(x^n, \varphi(f_0(x^n))) + \lambda \ell(f_0(x^n)) \\
&\leq \rho_n(x^n, \varphi(s)) + \ell(s) \\
&= \rho_n(x^n, f(x^n, z^m)) + \ell(f(x^n, z^m)).
\end{aligned}$$

Taking expectations, we see that $L_\theta(C_0^n, \lambda) \leq L_\theta(C^{n,m}, \lambda)$ for all $\theta \in \Lambda$, which proves that

$$\inf_{C^n} L_\theta(C^n, \lambda) \leq \inf_{m \geq 0} \inf_{C^{n,m}} L_\theta(C^{n,m}, \lambda), \qquad \forall \theta \in \Lambda.$$

The infimum of $L_\theta(C^n, \lambda)$ over all zero-memory variable-rate quantizers $C^n$ with block length $n$ is the *operational $n$th-order distortion-rate Lagrangian* $\widehat{L}_\theta^n(\lambda)$ [20]. Because each $P_\theta$ is ergodic, $\widehat{L}_\theta^n(\lambda)$ converges to the *distortion-rate Lagrangian*

$$L_\theta(\lambda) \triangleq \min_R \left( D_\theta(R) + \lambda R \right),$$

where $D_\theta(R)$ is the Shannon distortion-rate function of $P_\theta$ (see Lemma 2 in the Appendix to Chou, Effros and Gray [14]). Thus, our scheme is universal not only in the $n$th-order sense of (7), but also in the distortion-rate sense, i.e.,

$$L_\theta(C_*^{n,m_n}, \lambda) - L_\theta(\lambda) \to 0, \qquad \text{as } n \to \infty$$

for every $\theta \in \Lambda$. Thus, in the terminology of [14], our scheme is *weakly minimax universal* for $\{P_\theta : \theta \in \Lambda\}$.

## IV. PROOF OF THEOREM 3.1

### A. *The main idea*

In this section, we describe the main idea behind the proof and fix some notation. We have already seen that it suffices to construct a universal scheme that can compete with all *zero-memory* variable-rate quantizers. That is, it suffices to show that there exists a sequence $\{C_*^{n,m_n}\}$ of codes, such that

$$L_\theta(C_*^{n,m_n}, \lambda) - \widehat{L}_\theta^n(\lambda) = O\left( \sqrt{\frac{\mathsf{V}(\mathcal{A}_n)\log n}{n}} \right), \forall \theta \in \Lambda. \tag{8}$$

This is what we shall prove.

We assume throughout that the "true" source is $P_{\theta_0}$ for some $\theta_0 \in \Lambda$. Our code operates as follows. Suppose that:

- Both the encoder and the decoder have access to a countably infinite "database" $c = \{\theta(i)\}_{i \in \mathbb{N}}$, where each $\theta(i) \in \Lambda$. Using Elias' universal representation of the integers [36], we can associate to each $\theta(i)$ a unique binary string $s(i)$ with $\ell(s(i)) = \log i + O(\log \log i)$ bits.

- A sequence $\{\delta_n\}$ of positive reals is given, such that

$$\delta_n \to 0, \sqrt{n}\delta_n \to 0, \qquad \text{as } n \to \infty$$

(we shall specify the sequence $\{\delta_n\}$ later in the proof).

- For each $n \in \mathbb{N}$ and each $\theta \in \Lambda$, there exists a zero-memory $n$-block code $C_\theta^n = (f_\theta^n, \varphi_\theta^n)$ that achieves (or comes arbitrarily close to) the $n$th-order Lagrangian optimum for $P_\theta$: $L_\theta(C_\theta^n, \lambda) = \widehat{L}_\theta^n(\lambda)$.

Fix the block length $n$. Because the source is stationary, it suffices to describe the mapping of $X_1^n$ into $\widehat{X}_1^n$. The encoding is done as follows:

1) The encoder estimates $P_{\theta_0}^n$ from the $m_n$-block $X_{-m_n+1}^0$ as $P_{\widetilde{\theta}}^n$, where $\widetilde{\theta} = \widetilde{\theta}(X_{-m_n+1}^0)$.

2) The encoder then computes the *waiting time*

$$T_n \triangleq \inf \left\{ i \geq 1 : d_n\left( \theta(i), \widetilde{\theta}(X_{-m_n+1}^0) \right) \leq \sqrt{n}\delta_n \right\},$$

with the standard convention that the infimum of the empty set is equal to $+\infty$. That is, the encoder looks through the database $c$ and finds the first $\theta(i)$, such that the $n$-dimensional distribution $P_{\theta(i)}^n$ is in the variational ball of radius $\sqrt{n}\delta_n$ around $P_{\widetilde{\theta}}^n$.

3) If $T_n < +\infty$, the encoder sets $\widehat{\theta} = \theta(i)$; otherwise, the encoder sets $\widehat{\theta} = \theta_d$, where $\theta_d \in \Lambda$ is some default parameter vector, say, $\theta(1)$.

4) The binary description of $X_1^n$ is a concatenation of the following three binary strings: (i) a 1-bit flag $b$ to tell whether $T_n$ is finite ($b = 0$) or infinite ($b = 1$); (ii) a binary string $s_1$ which is equal to $s(T_n)$ if $T_n < +\infty$ or to an empty string if $T_n = +\infty$; (iii) $s_2 = f_{\widehat{\theta}}(X_1^n)$. The string $\widetilde{s} = bs_1$ is called the *first-stage description*, while $s_2$ is called the *second-stage description*.

The decoder receives $bs_1s_2$, determines $\widehat{\theta}$ from $\widetilde{s}$, and produces the reproduction $\widehat{X}_1^n = \varphi_{\widehat{\theta}}(s_2)$. Note that when $b = 0$ (which, as we shall show, will happen eventually almost surely), $P_{\widehat{\theta}}^n$ lies in the variational ball of radius $\sqrt{n}\delta_n$ around the estimated source $P_{\widetilde{\theta}}^n$. If the latter is a good estimate of $P_{\theta_0}^n$, i.e., $d_n(\theta_0, \widetilde{\theta}) \to 0$ as $n \to \infty$ a.s., then the estimate of the true source computed by the decoder is only slightly worse. Furthermore, as we shall show, the almost-sure convergence of $d_n(\theta_0, \widehat{\theta})$ to zero as $n \to \infty$ implies that the Lagrangian performance of $C_{\widehat{\theta}}^n$ on $P_{\theta_0}$ is close to the optimum $L_{\theta_0}(C_{\theta_0}^n, \lambda) \equiv \widehat{L}_{\theta_0}^n(\lambda)$.

Formally, the code $C_*^{n,m_n}$ is comprised by the following maps:

- the *parameter estimator* $\widetilde{\theta} : \mathcal{X}^{m_n} \to \Lambda$;
- the *parameter encoder* $\widetilde{g} : \Lambda \to \widetilde{\mathcal{S}}$, where $\widetilde{\mathcal{S}} \equiv \{0s(i)\}_{i \in \mathbb{N}} \cup \{1\}$;
- the *parameter decoder* $\widetilde{\psi} : \widetilde{\mathcal{S}} \to \Lambda$.

Let $\widetilde{f}$ denote the composition $\widetilde{g} \circ \widetilde{\theta}$ of the parameter estimator and the parameter encoder, which we refer to as the *first-stage encoder*, and let $\widehat{\theta}$ denote the composition $\widetilde{\psi} \circ \widetilde{f}$ of the parameter decoder and the first-stage encoder. The decoder $\widetilde{\psi}$ is the *first-stage decoder*. The collection $\{C_\theta^n : \theta \in \Lambda\}$ defines the *second-stage codes*. The encoder $f_* : \mathcal{X}^n \times \mathcal{X}^{m_n} \to \widetilde{\mathcal{S}} \times \mathcal{S}$

and the decoder $\varphi_* : \widetilde{S} \times S \to \widehat{\mathcal{X}}^n$ of $C_*^{n,m_n}$ are defined as

$$f_*(X_1^n, X_{-m_n+1}^0) \triangleq \widetilde{f}(X_{-m_n+1}^0) f_{\widehat{\theta}(X_{-m_n+1}^0)}(X_1^n)$$

and

$$\varphi_*(\widetilde{s}s) \triangleq \varphi_{\widetilde{\psi}(\widetilde{s})}(s), \qquad s \in S, \widetilde{s} \in \widetilde{S}$$

respectively. To assess the performance of $C_*^{n,m_n}$, consider the function

$$g(X_1^n, X_{-m_n+1}^0) \triangleq \rho_n\Big(X_1^n, C_{\widehat{\theta}(X_{-m_n+1}^0)}^n(X_1^n)\Big)$$
$$+\lambda\Big[\ell_n\Big(f_{\widehat{\theta}(X_{-m_n+1}^0)}(X_1^n)\Big) + \ell_n\Big(\widetilde{f}(X_{-m_n+1}^0)\Big)\Big].$$

The expectation $\mathbb{E}_{\theta_0}\big\{g(X_1^n, X_{-m_n+1}^0)\big\}$ of $g$ with respect to $P_{\theta_0}$ is precisely the Lagrangian performance of $C_*^{n,m_n}$, at Lagrange multiplier $\lambda$, on the source $P_{\theta_0}$. We consider separately the contributions of the first-stage and the second-stage codes. Define another function $h : \mathcal{X}^n \times \mathcal{X}^{m_n} \to \mathbb{R}^+$ by

$$h(X_1^n, X_{-m_n+1}^0) \triangleq \rho_n\Big(X_1^n, C_{\widehat{\theta}(X_{-m_n+1}^0)}^n(X_1^n)\Big)$$
$$+\lambda\ell_n\Big(f_{\widehat{\theta}(X_{-m_n+1}^0)}(X_1^n)\Big),$$

so that $\mathbb{E}_{\theta_0}\Big\{h(X_1^n, X_{-m_n+1}^0)\Big|X_{-m_n+1}^0\Big\}$ is the (random) Lagrangian performance of the code $C_{\widehat{\theta}(X_{-m_n+1}^0)}^n$ on $P_{\theta_0}$. Hence,

$$g(X_1^n, X_{-m_n+1}^0) = h(X_1^n, X_{-m_n+1}^0) + \lambda\ell_n\Big(\widetilde{f}(X_{-m+1}^0)\Big),$$

so, taking expectations, we get

$$L_{\theta_0}(C_*^{n,m_n}, \lambda) = \mathbb{E}_{\theta_0}\big\{h(X_1^n, X_{-m_n+1}^0)\big\}$$
$$+\lambda\,\mathbb{E}_{\theta_0}\Big\{\ell_n\Big(\widetilde{f}(X_{-m_n+1}^0)\Big)\Big\}. \tag{9}$$

Our goal is to show that the first term in Eq. (9) converges to the $n$th-order optimum $\widehat{L}_{\theta_0}^n(\lambda)$, and that the second term is $o(1)$.

The proof itself is organized as follows. First we motivate the choice of the memory lengths $m_n$ in Section IV-B. Then we indicate how to select the database $C$ (Section IV-C) and how to implement the parameter estimator $\widetilde{\theta}$ (Section IV-D) and the parameter encoder/decoder pair $(\widetilde{g}, \widetilde{\psi})$ (Section IV-E). The proof is concluded by estimating the Lagrangian performance of the resulting code (Section IV-F) and the fidelity of the source identification at the decoder (Section IV-G). In the following, (in)equalities involving the relevant random variables are assumed to hold for all realizations and not just a.s., unless specified otherwise.

### B. The memory length

Let $l_n = \lceil n^{(2+\eta)/r} \rceil$, where $r$ is the common decay exponent of the $\beta$-mixing coefficients $\beta_\theta(k)$ in Condition 1, and let $m_n = n(n + l_n)$. We divide the $m_n$-block $X_{-m_n+1}^0$ into $n$ blocks $Z_1, \ldots, Z_n$ of length $n$ interleaved by $n$ blocks $Y_1, \ldots, Y_n$ of length $l_n$ (see Figure 1). The parameter estimator $\widetilde{\theta}$, although defined as acting on the entire $X_{-m_n+1}^0$, effectively will make use only of $Z^n = (Z_1, \ldots, Z_n)$. The $Z_j$'s are each distributed according to $P_{\theta_0}^n$, but they are not

independent. Thus, the set

$$\mathcal{M} = \bigcup_{j=1}^n \{(j-1)(n + l_n) + 1 \le i \le j(n + l_n) - l_n\}$$

is the effective memory of $C_*^{n,m_n}$, and the effective memory length is $n^2$.

Let $Q^{(n)}$ denote the marginal distribution of $Z^n$, and let $\widetilde{Q}^{(n)}$ denote the product of $n$ copies of $P_{\theta_0}^n$. We now show that we can approximate $Q^{(n)}$ by $\widetilde{Q}^{(n)}$ in variational distance, increasingly finely with $n$. Note that both $Q^{(n)}$ and $\widetilde{Q}^{(n)}$ are defined on the $\sigma$-algebra $\mathcal{F}^{(n)}$, generated by all $X_i$ except those in $Y_1, \ldots, Y_n$, so that $d(Q^{(n)}, \widetilde{Q}^{(n)}) = d(Q^{(n)}, \widetilde{Q}^{(n)}; \mathcal{F}^{(n)})$. Therefore, using induction and the definition of the $\beta$-mixing coefficient (cf. Section II-A), we have

$$d(Q^{(n)}, \widetilde{Q}^{(n)}) \le (n-1)\beta_{\theta_0}(l_n) = O(1/n^{1+\eta}),$$

where the last equality follows from Condition 1 and from our choice of $l_n$. This in turn implies the following useful fact (see also Lemma 4.1 of Yu [29]), which we shall heavily use in the proof: for any measurable function $\sigma : \mathcal{X}^{n^2} \to [0, M]$ with $M < \infty$,

$$\Big|\mathbb{E}_{Q^{(n)}}\big\{\sigma(Z^n)\big\} - \mathbb{E}_{\widetilde{Q}^{(n)}}\big\{\sigma(Z^n)\big\}\Big| \le M(n-1)\beta_{\theta_0}(l_n)$$
$$= O(1/n^{1+\eta}), \tag{10}$$

where the constant hidden in the $O(\cdot)$ notation depends on $M$ and on $\theta_0$.

### C. Construction of the database

The database, or the first-stage codebook, $c$ is constructed by random selection. Let $W$ be a probability distribution on $\Lambda$ which is absolutely continuous with respect to the Lebesgue measure and has an everywhere positive and continuous density $w(\theta)$. Let $C = \{\Theta(i)\}_{i \in \mathbb{N}}$ be a collection of independent random vectors taking values in $\Lambda$, each generated according to $W$ independently of $X$. We use $\mathbb{W}$ to denote the process distribution of $C$.

Note that the first-stage codebook is countably infinite, which means that, in principle, both the encoder and the decoder must have unbounded memory in order to store it. This difficulty can be circumvented by using synchronized random number generators at the encoder and at the decoder, so that the entries of $C$ can be generated as needed. Thus, by construction, the encoder will generate $T_n$ samples (where $T_n$ is the waiting time) and then communicate (a binary encoding of) $T_n$ to the decoder. Since the decoder's random number generator is synchronized with that of the encoder's, the decoder will be able to recover the required entry of $C$.

### D. Parameter estimation

The parameter estimator $\widetilde{\theta} : \mathcal{X}^{m_n} \to \Lambda$ is constructed as follows. Because the source $X$ is stationary, it suffices to describe the action of $\widetilde{\theta}$ on $X_{-m_n+1}^0$. In the notation of Section IV-A, let $P_{Z^n}$ be the empirical distribution of
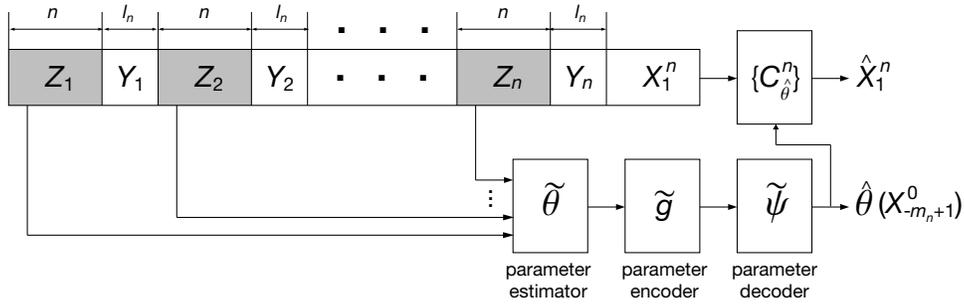
Fig. 1. The structure of the code $C_*^{n,m_n}$. The shaded blocks are those used for estimating the source parameters.

$Z^n = (Z_1, \ldots, Z_n)$. For every $\theta \in \Lambda$, define

$$U_\theta(Z^n) \triangleq \sup_{A \in \mathcal{A}_n} |P_\theta^n(A) - P_{Z^n}(A)|$$

$$\equiv \sup_{A \in \mathcal{A}_n} \left| \int_A p_\theta(x^n) d\mu(x^n) - P_{Z^n}(A) \right|,$$

where $\mathcal{A}_n$ is the Yatracos class defined by the $n$th-order densities $\{p_\theta^n : \theta \in \Lambda\}$ (see Section III). Finally, define $\widetilde{\theta}(X_{-m+1}^0)$ as any $\theta^* \in \Lambda$ satisfying

$$U_{\theta^*}(Z^n) < \inf_{\theta \in \Lambda} U_\theta(Z^n) + \frac{1}{n},$$

where the extra $1/n$ term is there to ensure that at least one such $\theta^*$ exists. This is the so-called *minimum-distance (MD) density estimator* of Devroye and Lugosi [34], [35] (see also Devroye and Györfi [37]), adapted to the dependent-process setting of the present paper. The key property of the MD estimator is that

$$d_n\left(\widetilde{\theta}(X_{m_n+1}^0), \theta_0\right) \leq 4U_{\theta_0}(Z^n) + \frac{3}{n} \qquad (11)$$

(see, e.g., Theorem 5.1 of Devroye and Györfi [37]). This holds regardless of whether the samples $Z_1, \ldots, Z_n$ are independent or not.

### E. Encoding and decoding of parameter estimates

Next we construct the parameter encoder-decoder pair $(\widetilde{g}, \widetilde{\psi})$. Given a $\theta \in \Lambda$, define the *waiting time*

$$T_n(\theta) \triangleq \inf\{i \geq 1 : d_n(\theta, \Theta(i)) \leq \sqrt{n}\delta_n\},$$

with the standard convention that the infimum of the empty set is equal to $+\infty$. That is, given a $\theta \in \Lambda$, the parameter encoder looks through the codebook $C$ and finds the position of the first $\Theta(i)$ such that the variational distance between the $n$th-order distributions $P_\theta^n$ and $P_{\Theta(i)}^n$ is at most $\sqrt{n}\delta_n$. If no such $\Theta(i)$ is found, the encoder sets $T_n = +\infty$. We then define the maps $\widetilde{g}$ and $\widetilde{\psi}$ by

$$\widetilde{g}(\theta) = \begin{cases} 0s(T_n), & \text{if } T_n < \infty \\ 1, & \text{if } T_n = \infty \end{cases}$$

and

$$\widetilde{\psi}(0s(i)) = \Theta(i), \qquad \widetilde{\psi}(1) = \theta(1)$$

respectively. Thus, $\widetilde{S} = \{0s(i)\} \cup \{1\}$, and the bound

$$\ell(\widetilde{g}(\theta)) \leq \log T_n + O(\log \log T_n) \qquad (12)$$

holds for every $\theta \in \Lambda$, regardless of whether $T_n$ is finite or infinite.

### F. Performance of the code

Given the random codebook $C$, the expected Lagrangian performance of our code on the source $P_{\theta_0}$, is

$$L_{\theta_0}(C_*^{n,m_n}, \lambda) = \mathbb{E}_{\theta_0}\left\{ g\left(X_1^n, X_{-m_n+1}^0\right) \right\}$$

$$= \mathbb{E}_{\theta_0}\left\{ h\left(X_1^n, X_{-m_n+1}^0\right) \right\}$$

$$+ \lambda \mathbb{E}_{\theta_0}\left\{ \ell_n\left(\widetilde{f}\left(X_{-m_n+1}^0\right)\right) \right\}. \qquad (13)$$

We now upper-bound the two terms in (13). We start with the second term.

We need to bound the expectation of the waiting time $T_n = T_n(\widetilde{\theta}(X_{-m_n+1}^0))$. Our strategy borrows some elements from the paper of Kontoyiannis and Zhang [38]. Consider the probability

$$q_n \triangleq W\left(d_n\big(\Theta, \widetilde{\theta}(X_{-m_n+1}^0)\big) \leq \sqrt{n}\delta_n\right),$$

which is a random function of $X_{-m_n+1}^0$. From Condition 2, it follows for $n$ sufficiently large that

$$q_n \geq W\left(\|\Theta - \widetilde{\theta}(X_{-m_n+1}^0)\| \leq \delta_n/c_{\widetilde{\theta}}\right),$$

where $\widetilde{\theta} \equiv \widetilde{\theta}(X_{-m_n+1}^0)$. Because the density $w(\theta)$ is everywhere positive, the latter probability is strictly positive for almost all $X_{-m_n+1}^0$, and so $q_n > 0$ eventually almost surely. Thus, the waiting times $T_n$ will be finite eventually almost surely (with respect to both the source $X$ and the first-stage codebook $C$). Now, if $q_n > 0$, then, conditioned on $X_{-m_n+1}^0 = x_{-m_n+1}^0$, the waiting time $T_n$ is a geometric random variable with parameter $q_n$, and it is not hard to show (see, e.g., Lemma 3 of Kontoyiannis and Zhang [38]) that for any $\epsilon > 0$

$$\mathbb{P}\left(\log[(T_n - 1)q_n] \geq \epsilon \Big| X_{-m_n+1}^0 = x_{-m_n+1}^0\right) \leq e^{-2^\epsilon}.$$

Setting $\epsilon = \log(2\log n)$, we have, for almost all $X_{-m_n+1}^n$, that

$$\mathbb{P}\left(\log[(T_n - 1)q_n] \geq \log(2\log n)\Big| X_{-m_n+1}^0 = x_{-m_n+1}^0\right)$$

$$\leq e^{-2\log n} \leq n^{-2}.$$

Then, by the Borel–Cantelli lemma,

$$\log(T_n q_n) \leq \log\log n + 2$$

eventually almost surely, so that

$$\mathbb{E}_{\theta_0}\{\log T_n\} \le \log\log n + 2 - \mathbb{E}_{\theta_0}\{\log q_n\} \qquad (14)$$

for almost every realization of the random codebook $\boldsymbol{C}$ and for sufficiently large $n$. We now obtain an asymptotic lower bound on $\mathbb{E}_{\theta_0}\{\log q_n\}$. Define the events

$$
\begin{aligned}
F_n &\triangleq \left\{ (d_n\big(\widetilde{\theta}(X^0_{-m_n+1}), \theta_0) \le \sqrt{n}\delta_n/2 \right\}, \\
G_n &\triangleq \left\{ d_n(\Theta, \theta_0) \le \sqrt{n}\delta_n/2 \right\}, \\
H_n &\triangleq \left\{ \|\Theta - \theta_0\| \le \delta_n/2c_{\theta_0} \right\}.
\end{aligned}
$$

Then by the triangle inequality we have

$$F_n \text{ and } G_n \implies d_n\big(\Theta, \widetilde{\theta}(X^0_{-m_n+1})\big) \le \sqrt{n}\delta_n,$$

and, for $n$ sufficiently large, we can write

$$
\begin{aligned}
q_n &\overset{(a)}{\ge} W(G_n)P_{\theta_0}(F_n) \\
&\overset{(b)}{=} W(G_n)Q^{(n)}(F_n) \\
&\overset{(c)}{\ge} W(H_n)Q^{(n)}(F_n),
\end{aligned}
$$

where (a) follows from the independence of $\boldsymbol{X}$ and $\boldsymbol{C}$, (b) follows from the fact that the parameter estimator $\widetilde{\theta}(X^0_{-m_n+1})$ depends only on $Z^n$, and (c) follows from Condition 2 and the fact that $\delta_n \to 0$. Since the density $w$ is everywhere positive and continuous at $\theta_0$, $w(\theta) \ge w(\theta_0)/2$ for all $\theta \in H_n$ for $n$ sufficiently large, so

$$W(H_n) = \int_{H_n} w(\theta)d\theta \ge \frac{1}{2}w(\theta_0)v_k\left(\frac{\delta_n}{2c_{\theta_0}}\right)^k, \qquad (15)$$

where $v_k$ is the volume of the unit sphere in $\mathbb{R}^k$. Next, the fact that the minimum-distance estimate $\widetilde{\theta}(X^0_{-m_n+1})$ depends only on $Z^n$ implies that the event $F_n$ belongs to the $\sigma$-algebra $\mathcal{F}^{(n)}$, and from (10) we get

$$Q^{(n)}(F_n) \ge \widetilde{Q}^{(n)}(F_n) - O(1/n^{1+\eta}). \qquad (16)$$

Under $\widetilde{Q}^{(n)}$, the $n$-blocks $Z_1, \ldots, Z_n$ are i.i.d. according to $P^n_{\theta_0}$, and we can invoke the Vapnik–Chervonenkis machinery to lower-bound $\widetilde{Q}^{(n)}(F_n)$. In the notation of Sec. IV-D, define the event

$$I_n \triangleq \left\{ 4U_{\theta_0}(Z^n) + \frac{3}{n} \le \frac{\sqrt{n}\delta_n}{2} \right\}.$$

Then $I_n$ implies $F_n$ by (11), and

$$\widetilde{Q}^{(n)}(F_n^c) \le \widetilde{Q}^{(n)}(I_n^c) \le 8n^{\mathsf{V}(\mathcal{A}_n)}e^{-n(\sqrt{n}\delta_n - 6/n)^2/2048}, \qquad (17)$$

where the second bound is by the Vapnik–Chervonenkis inequality (3) of Lemma 2.1. Combining the bounds (16) and (17) and using Condition 1, we obtain

$$P^m_\theta(F_n) \ge 1 - 8n^{\mathsf{V}(\mathcal{A}_n)}e^{-n(\sqrt{n}\delta_n - 6/n)^2/2048} - O(1/n^{1+\eta}) \qquad (18)$$

Now, if we choose

$$\delta_n = \frac{\sqrt{2048(\mathsf{V}(\mathcal{A}_n)+1)\ln n}}{n} + \frac{6}{n^{3/2}},$$

then the right-hand side of (18) can be further lower-bounded

by $1 - O(1/n)$. Combining this with (15), taking logarithms, and then taking expectations, we obtain

$$
\begin{aligned}
\mathbb{E}_{\theta_0}\{\log q_n\} \\
\ge\ & \log(1 - O(1/n)) + k\log\delta_n + 2^{c(k,\theta_0)} \\
=\ & \log(1 - O(1/n)) \\
& + k\log\left[\sqrt{2048(\mathsf{V}(\mathcal{A}_n)+1)n\ln n} + 6\right] \\
& + \frac{3k}{2}\log\frac{1}{n} + c(k,\theta_0) \\
\ge\ & \log(1 - O(1/n)) + \frac{3k}{2}\log\frac{1}{n} + c(k,\theta_0),
\end{aligned}
$$

where $c(k,\theta_0)$ is a constant that depends only on $k$ and $\theta_0$. Using this and (14), we get that

$$\mathbb{E}_{\theta_0}\{\log T_n\} \le \log\log n + O(\log n)$$

for $\mathbb{W}$-almost every realization of the random codebook $\boldsymbol{C}$, for $n$ sufficiently large. Together with (12), this implies that

$$
\begin{aligned}
\mathbb{E}_{\theta_0}\left\{\ell_n\left(\widetilde{f}\left(X^0_{-m_n+1}\right)\right)\right\} \\
=\ O\left(\frac{\log n}{n}\right) + O\left(\frac{\log\log n}{n}\right) + \frac{3}{n} + o(1)
\end{aligned}
$$

for $\mathbb{W}$-almost all realizations of the first-stage codebook.

We now turn to the first term in (13). Recall that, for each $\theta \in \Lambda$, the code $C^n_\theta$ is $n$th-order optimal for $P_\theta$. Using this fact together with the boundedness of the distortion measure $\rho$, we can invoke Lemma A.3 in the Appendix and assume without loss of generality that each $C^n_\theta$ has a finite codebook (of size not exceeding $2^{n\rho_{\max}/\lambda}$), and each codevector can be described by a binary string of no more than $2n\rho_{\max}/\lambda$ bits. Hence, $h(X^n_1, X^0_{-m_n+1}) \le 3\rho_{\max}$. Let $P^-$ and $P^+$ be the marginal distributions of $P_{\theta_0}$ on $\sigma(X^0_{-\infty})$ and $\sigma(X^\infty_1)$, respectively. Note that $h(X^n, X^0_{-m_n+1})$ does not depend on $X^0_{-l_n+1}$. This, together with Condition 1 and the choice of $l_n$, implies that

$$
\begin{aligned}
\mathbb{E}_{\theta_0}\left\{h\left(X^n_1, X^0_{-m_n+1}\right)\right\} \\
\le \mathbb{E}_{P^-\times P^+}\left\{h\left(X^n_1, X^0_{-m_n+1}\right)\right\} + \beta_{\theta_0}(l_n) \\
= \mathbb{E}_{P^-\times P^+}\left\{h\left(X^n_1, X^0_{-m_n+1}\right)\right\} + O(1/n^{2+\eta}).
\end{aligned}
$$

Furthermore,

$$
\begin{aligned}
\mathbb{E}_{P^-\times P^+}&\left\{h\left(X^n_1, X^0_{-m_n+1}\right)\right\} \\
&= \int_{\mathcal{X}^n\times\mathcal{X}^{m_n}} h(x^n, z^{m_n})dP_{\theta_0}(x^n)dP_{\theta_0}(z^{m_n}) \\
&\overset{(a)}{=} \int_{\mathcal{X}^{m_n}} \mathbb{E}_{\theta_0}\{h(X^n_1, z^{m_n})\}dP_{\theta_0}(z^{m_n}) \\
&\overset{(b)}{=} \mathbb{E}_{\theta_0}\left\{L_{\theta_0}\left(C^n_{\widehat{\theta}(X^0_{-m_n+1})}, \lambda\right)\right\},
\end{aligned}
$$

where (a) follows by Fubini's theorem and the boundedness of $h$, while (b) follows from the definition of $h$. The Lagrangian performance of the code $C^n_{\widehat{\theta}}$, where $\widehat{\theta} = \widehat{\theta}(X^0_{-m_n+1})$, can be

further bounded as

$$L_{\theta_0}\left(C_{\widehat{\theta}}^n, \lambda\right)$$

$$\overset{(a)}{\leq} L_{\widehat{\theta}}\left(C_{\widehat{\theta}}^n, \lambda\right) + 3\rho_{\max} d_n\left(\widehat{\theta}(X^0_{-m_n+1}), \theta_0\right)$$

$$\overset{(b)}{=} \widehat{L}_{\widehat{\theta}}^n(\lambda) + 3\rho_{\max} d_n\left(\widehat{\theta}(X^0_{-m_n+1}), \theta_0\right)$$

$$\overset{(c)}{\leq} \widehat{L}_{\theta_0}^n(\lambda) + 4\rho_{\max} d_n\left(\widehat{\theta}(X^0_{-m_n+1}), \theta_0\right)$$

$$\overset{(d)}{\leq} \widehat{L}_{\theta_0}^n(\lambda) + 4\rho_{\max}\Big[d_n\left(\widehat{\theta}(X^0_{-m_n+1}), \widetilde{\theta}(X^0_{-m_n+1})\right)$$
$$+ d_n\left(\widetilde{\theta}(X^0_{-m_n+1}), \theta_0\right)\Big],$$

where (a) follows from Lemma A.3 in the Appendix, (b) follows from the $n$th-order optimality of $C_{\widehat{\theta}}^n$ for $P_{\widehat{\theta}}^n$, (c) follows, overbounding slightly, from the Lagrangian mismatch bound of Lemma A.2 in the Appendix, and (d) follows from the triangle inequality. Taking expectations, we obtain

$$\mathbb{E}_{\theta_0}\left\{L_{\theta_0}\left(C_{\widehat{\theta}(X^0_{-m_n+1})}^n, \lambda\right)\right\} \leq \widehat{L}_{\theta_0}^n(\lambda)$$
$$+ 4\rho_{\max} \cdot \mathbb{E}_{\theta_0}\left\{d_n\left(\widehat{\theta}(X^0_{-m_n+1}), \widetilde{\theta}(X^0_{-m_n+1})\right)\right.$$
$$\left. + d_n\left(\widetilde{\theta}(X^0_{-m_n+1}), \theta_0\right)\right\}. \tag{19}$$

The second $d_n(\cdot, \cdot)$ term in (19) can be interpreted as the estimation error due to estimating $P_{\theta_0}^n$ by $P_{\widehat{\theta}}^n$, while the first $d_n(\cdot, \cdot)$ is the approximation error due to quantization of the parameter estimate $\widetilde{\theta}$. We examine the estimation error first. Using (11), we can write

$$\mathbb{E}_{\theta_0}\left\{d\left(P_{\theta^*(X^0_{-m+1})}^n, P_{\theta_0}^n\right)\right\} \leq 4\,\mathbb{E}_{\theta_0}\{U_{\theta_0}(Z^n)\} + \frac{3}{n}. \tag{20}$$

Now, each $Z_j$ is distributed according to $P_{\theta_0}^n$, and we can approximate the expectation of $U_{\theta_0}(Z^n)$ with respect to $Q^{(n)}$ by the expectation of $U_{\theta_0}(Z^n)$ with respect to the product measure $\widetilde{Q}^{(n)}$:

$$\mathbb{E}_{Q^{(n)}}\{U_{\theta_0}(Z^n)\} \leq \mathbb{E}_{\widetilde{Q}^{(n)}}\{U_{\theta_0}(Z^n)\} + (n-1)\beta_\theta(l_n)$$
$$\leq c\sqrt{\frac{\mathsf{V}(\mathcal{A}_n)\log n}{n}} + O\left(\frac{1}{n^{1+\eta}}\right)$$
$$= O\left(\sqrt{\frac{\mathsf{V}(\mathcal{A}_n)\log n}{n}}\right),$$

where the second estimate follows from the Vapnik–Chervonenkis inequality (4) and from the choice of $l_n$. This, together with (20), yields

$$\mathbb{E}_{\theta_0}\left\{d\left(P_{\theta^*(X^0_{-m+1})}^n, P_{\theta_0}^n\right)\right\} = O\left(\sqrt{\frac{\mathsf{V}(\mathcal{A}_n)\log n}{n}}\right). \tag{21}$$

As for the first $d_n(\cdot, \cdot)$ term in (19), we have, by construction of the first-stage encoder, that

$$d_n\left(\widehat{\theta}(X^0_{-m_n+1}), \widetilde{\theta}(X^0_{-m_n+1})\right)$$
$$\leq \sqrt{n}\delta_n = O\left(\sqrt{\frac{\mathsf{V}(\mathcal{A}_n)\log n}{n}}\right) \tag{22}$$

eventually almost surely, so the corresponding expectation is $O\big(\sqrt{\mathsf{V}(\mathcal{A}_n)\log n/n}\big)$ as well. Summing the estimates (21) and (22), we obtain

$$\mathbb{E}_{\theta_0}\left\{h\left(X_1^n, X^0_{-m_n+1}\right)\right\} = \widehat{L}_{\theta_0}^n(\lambda) + O\left(\sqrt{\frac{\mathsf{V}(\mathcal{A}_n)\log n}{n}}\right).$$

Finally, putting everything together, we see that, eventually,

$$L_{\theta_0}\left(C_*^{n,m_n}\right) = \widehat{L}_{\theta_0}^n(\lambda) + O\left(\sqrt{\frac{\mathsf{V}(\mathcal{A}_n)\log n}{n}}\right)$$
$$+ \lambda\Big[O\Big(\frac{\log n}{n}\Big) + O\Big(\frac{\log\log n}{n}\Big) + \frac{3}{n} + o(1)\Big] \tag{23}$$

for $\mathbb{W}$-almost every realization of the first-stage codebook $C$. This proves (8), and hence (7).

### G. Identification of the active source

We have seen that the expected variational distance $\mathbb{E}_{\theta_0}\left\{d_n\left(\theta_0, \widehat{\theta}(X^0_{-m_n+1})\right)\right\}$ between the $n$-dimensional marginals of the true source $P_{\theta_0}$ and the estimated source $P_{\widehat{\theta}(X^0_{-m_n+1})}$ converges to zero as $\sqrt{\mathsf{V}(\mathcal{A}_n)\log n/n}$. We wish to show that this convergence also holds eventually with probability one, i.e.,

$$d_n(\theta_0, \widehat{\theta}(X^0_{-m_n+1})) = O\left(\sqrt{\frac{\mathsf{V}(\mathcal{A}_n)\log n}{n}}\right) \tag{24}$$

$P_{\theta_0}$-almost surely.

Given an $\epsilon > 0$, we have by the triangle inequality that $d_n(\theta_0, \widehat{\theta}(X^0_{-m_n+1})) > \epsilon$ implies

$$d_n\left(\theta_0, \widetilde{\theta}(X^0_{-m_n+1})\right) + d_n\left(\widetilde{\theta}(X^0_{-m_n+1}), \widehat{\theta}(X^0_{-m_n+1})\right) > \epsilon,$$

where $\widetilde{\theta}(X^0_{-m_n+1})$ is the minimum-distance estimate of $P_{\theta_0}^n$ from $X^0_{-m_n+1}$ (cf. Section IV-E). Recalling our construction of the first-stage encoder, we see that this further implies

$$d_n\left(\theta_0, \widetilde{\theta}(X^0_{-m_n+1})\right) > \epsilon - \sqrt{n}\delta_n.$$

Finally, using the property (11) of the minimum-distance estimator, we obtain that

$$d_n\left(\theta_0, \widehat{\theta}(X^0_{-m_n+1})\right) > \epsilon$$

implies

$$U_{\theta_0}(Z^n) > \frac{1}{4}\left(\epsilon - \sqrt{n}\delta_n - \frac{3}{n}\right).$$

Therefore,

$$Q^{(n)}\left\{d_n\left(\theta_0, \widehat{\theta}(X^0_{-m_n+1})\right) > \epsilon\right\}$$
$$\leq Q^{(n)}\left\{U_{\theta_0}(Z^n) > \frac{1}{4}\left(\epsilon - \sqrt{n}\delta_n - \frac{3}{n}\right)\right\}$$
$$\overset{(a)}{\leq} \widetilde{Q}^{(n)}\left\{U_{\theta_0}(Z^n) > \frac{1}{4}\left(\epsilon - \sqrt{n}\delta_n - \frac{3}{n}\right)\right\}$$
$$+ (n-1)\beta_{\theta_0}(l_n)$$
$$\overset{(b)}{\leq} 8n^{\mathsf{V}(\mathcal{A}_n)}\exp\left(-\frac{n(\epsilon - \sqrt{n}\delta_n - 3/n)^2}{512}\right)$$
$$+ (n-1)\beta_{\theta_0}(l_n), \tag{25}$$

where (a) follows, as before, from the definition of the $\beta$-mixing coefficient and (b) follows by the Vapnik–Chervonenkis inequality. Now, if we choose

$$\epsilon_n = \sqrt{\frac{512(\mathsf{V}(\mathcal{A}_n)\ln n + n\delta)}{n}} + \sqrt{n}\delta_n + \frac{3}{n}$$

for an arbitrary small $\delta > 0$, then (25) can be further upper-bounded by $8e^{-n\delta} + \sum_n n\beta_\theta(\ell_n)$, which, owing to Condition 1 and the choice $l_n = \lceil n^{(2+\eta)/r}\rceil$, is summable in $n$. Thus,

$$\sum_n Q^{(n)}\left\{d_n\left(\theta_\theta, \widehat{\theta}(X^0_{-m_n+1})\right) > \epsilon_n\right\} < \infty,$$

and we obtain (24) by the Borel–Cantelli lemma.

## V. EXAMPLES

### A. Stationary memoryless sources

As a basic check, let us see how Theorem 3.1 applies to stationary memoryless (i.i.d.) sources. Let $\mathcal{X} = \mathbb{R}$, and let $\{P_\theta : \theta \in \Lambda\}$ be the collection of all Gaussian i.i.d. processes, where

$$\Lambda = \{(m, \sigma) : m \in \mathbb{R}, 0 < \sigma < \infty\} \subset \mathbb{R}^2.$$

Then the $n$-dimensional marginal for a given $\theta = (m, \sigma)$ has the Gaussian density

$$p_\theta(x^n) = \frac{1}{(2\pi\sigma^2)^{n/2}} \prod_{i=1}^{n} e^{-(x_i - m)^2/2\sigma^2}$$

with respect to the Lebesgue measure. This class of sources trivially satisfies Condition 1 with $r = +\infty$, and it remains to check Conditions 2 and 3.

To check Condition 2, let us examine the normalized $n$th-order relative entropy between $P_\theta$ and $P_{\theta'}$, with $\theta = (m, \sigma)$ and $\theta' = (m', \sigma')$. Because the sources are i.i.d.,

$$D_n(P_\theta \| P_{\theta'}) = D(P_\theta^1 \| P_{\theta'}^1)$$
$$= \frac{1}{2}\left(\ln\left(\frac{\sigma}{\sigma'}\right)^2 + \left(\frac{\sigma'}{\sigma}\right)^2 + \frac{(m - m')^2}{\sigma'^2} - 1\right).$$

Applying the inequality $\ln x \le x - 1$ and some straightforward algebra, we get the bound

$$D_n(P_\theta \| P_{\theta'}) \le \left(\frac{\sigma + \sigma'}{\sigma}\right)^2 \frac{(\sigma - \sigma')^2}{2\sigma'^2} + \frac{(m - m')^2}{2\sigma'^2}$$
$$\le \left(1 + \frac{\sigma'}{\sigma}\right)^2 \frac{\|\theta - \theta'\|^2}{2\sigma'^2}.$$

Now fix a small $\delta \in (0, \sigma)$, and suppose that $\|\theta - \theta'\| < \delta$. Then $|\sigma - \sigma'| < \delta$, so we can further upper-bound $D_n(P_\theta \| P_{\theta'})$ by

$$D_n(P_\theta \| P_{\theta'}) \le \frac{9}{2(\sigma - \delta)^2}\|\theta - \theta'\|^2.$$

Thus, for a given $\theta \in \Lambda$, we see that

$$D_n(P_\theta \| P_{\theta'}) \le \frac{c_\theta^2}{2}\|\theta - \theta'\|^2$$

for all $\theta'$ in the open ball of radius $\delta$ around $\theta$, with $c_\theta \triangleq$ $3/(\sigma - \delta)$. Using Pinsker's inequality, we have

$$\frac{d_n(\theta, \theta')}{\sqrt{n}} \le \sqrt{2D_n(P_\theta \| P_{\theta'})} \le c_\theta \|\theta - \theta'\|$$

for all $n$. Thus, Condition 2 holds.

To check Condition 3, note that, for each $n$, the Yatracos class $\mathcal{A}_n$ consists of all sets of the form

$$\left\{x^n \in \mathbb{R}^n : \ln\sigma^2 - \ln\sigma'^2 + \frac{1}{\sigma^2}\sum_{i=1}^{n}(x_i - m)^2 \right.$$
$$\left. -\frac{1}{\sigma'^2}\sum_{i=1}^{n}(x_i - m')^2 > 0\right\} \qquad (26)$$

for all $m, m' \in \mathbb{R}; \sigma, \sigma' \in (0, \infty)$. Let $\alpha \triangleq \ln\sigma^2$ and $\alpha' \triangleq \ln\sigma'^2$. Then we can rewrite (26) as

$$\left\{x^n \in \mathbb{R}^n : \alpha - \alpha' + \frac{1}{\sigma^2}\sum_{i=1}^{n}(x_i - m)^2 \right.$$
$$\left. -\frac{1}{\sigma'^2}\sum_{i=1}^{n}(x_i - m')^2 > 0\right\}.$$

This is the set of all $x^n \in \mathbb{R}^n$ such that

$$\Pi(x^n, \alpha, \alpha', 1/\sigma^2, 1/\sigma'^2, m, m') > 0,$$

where $\Pi(x^n, \cdot)$ is a third-degree polynomial in the six parameters $(\alpha, \alpha', 1/\sigma^2, 1/\sigma'^2, m, m')$. It then follows from Lemma 2.2 that $\mathcal{A}_n$ is a VC class with $\mathsf{V}(\mathcal{A}_n) \le 12\log(12e)$. Therefore, Condition 3 holds as well.

### B. Autoregressive sources

Again, let $\mathcal{X} = \mathbb{R}$ and consider the case when $\boldsymbol{X}$ is a Gaussian autoregressive source of order $p$, i.e., it is the output of an autoregressive filter of order $p$ driven by white Gaussian noise. Then there exist $p$ real parameters $a_1, \ldots, a_p$ (the filter coefficients), such that

$$X_n = -\sum_{i=1}^{p} a_i X_{n-i} + Y_n, \qquad \forall n \in \mathbb{N}$$

where $\boldsymbol{Y} = \{Y_i\}_{i \in \mathbb{Z}}$ is an i.i.d. Gaussian process with zero mean and unit variance. Let $\Lambda \subset \mathbb{R}^p$ be the set of all $a_1, \ldots, a_p$, such that the roots of the polynomial $A(z) = \sum_{i=0}^{p} a_i z^i$, where $a_0 \equiv 1$, lie outside the unit circle in the complex plane. This ensures that $\boldsymbol{X}$ is a stationary process. We now proceed to check that Conditions 1–3 of Section III are satisfied.

The distribution of each $Y_i$ is absolutely continuous, and we can invoke the result of Mokkadem [31] to conclude that, for each $\theta \in \Lambda$, the process $\boldsymbol{X}$ is *geometrically mixing*, i.e., for every $\theta \in \Lambda$, there exists some $\gamma = \gamma(\theta) \in (0, 1)$, such that $\beta_\theta(k) = O(\gamma^k)$. Now, for any fixed $r > 0$, $\gamma^k \le k^{-r}$ for $k$ sufficiently large, so Condition 1 holds.

To check Condition 2, note that, for each $\theta \in \Lambda$, the Fisher information matrix $I_n(\theta)$ is independent of $n$ (see, e.g., Section 6 of Klein and Spreij [39]). Thus, the asymptotic Fisher information matrix $I(\theta) = \lim_{n \to \infty} I_n(\theta)$ exists and is nonsingular [39, Theorem 6.1], so, recalling the discussion in Section III, we conclude that Condition 2 holds also.

To verify Condition 3, consider the $n$-dimensional marginal $P_\theta(x^n)$, which has the Gaussian density

$$p_\theta(x^n) = \frac{1}{(2\pi \det R_n(\theta))^{n/2}} e^{-\frac{1}{2}x^{nT}R_n^{-1}(\theta)x^n},$$

where $R_n(\theta) \equiv \mathbb{E}_\theta\left\{(X_1^n)^T X_1^n\right\}$ is the $n$th-order autocorrelation matrix of $\mathbf{X}$. Thus, the Yatracos class $\mathcal{A}_n$ consists of sets of the form

$$A_{\theta,\theta'} = \left\{ x^n \in \mathbb{R}^n : \frac{n}{2}\ln\det R_n^{-1}(\theta) - \frac{1}{2}x^{nT}R_n^{-1}(\theta)x^n \right.$$
$$\left. > \frac{n}{2}\ln\det R_n^{-1}(\theta') - \frac{1}{2}x^{nT}R_n^{-1}(\theta')x^n \right\}$$

for all $\theta, \theta' \in \Lambda$. Now, for every $\theta \in \Lambda$, let $\bar{\theta} \triangleq (\theta, \ln\det R_n^{-1}(\theta))$. Since $\ln\det R_n^{-1}(\theta)$ is uniquely determined by $\theta$, we have $A_{\theta,\theta'} = A_{\bar{\theta},\bar{\theta}'}$ for all $\theta, \theta' \in \Lambda$. Using this fact, as well as the easily established fact that the entries of the inverse covariance matrix $R_n^{-1}(\theta)$ are second-degree polynomials in the filter coefficients $a_1, \ldots, a_p$, we see that, for each $x^n$, the condition $x^n \in A_{\theta,\theta'}$ can be expressed as $\Pi(x^n, \bar{\theta}) > 0$, where $\Pi(x^n, \cdot)$ is quadratic in the $2p+2$ real variables $\bar{\theta}_1, \ldots, \bar{\theta}_{p+1}, \bar{\theta}'_1, \ldots, \bar{\theta}'_{p+1}$. Thus, we can apply Lemma 2.2 to conclude that $\mathsf{V}(\mathcal{A}_n) \leq (4p+4)\log(8e)$. Therefore, Condition 3 is satisfied as well.

### C. Hidden Markov processes

A hidden Markov process (or a hidden Markov model, see, e.g., [40]) is a discrete-time bivariate random process $\{(S_i, X_i)\}$, where $\mathbf{S} = \{S_i\}$ is a homogeneous Markov chain and $\mathbf{X} = \{X_i\}$ is a sequence of random variables which are conditionally independent given $\mathbf{S}$, and the conditional distribution of $X_n$ is time-invariant and depends on $\mathbf{S}$ only through $S_n$. The Markov chain $\mathbf{S}$, also called the *regime*, is not available for observation. The observable component $\mathbf{X}$ is the source of interest. In information theory (see, e.g., [41] and references therein), a hidden Markov process is a discrete-time finite-state homogeneous Markov chain $\mathbf{S}$, observed through a discrete-time memoryless channel, so that $\mathbf{X} = \{X_i\}$ is the observation sequence at the output of the channel.

Let $M$ denote the number of states of $\mathbf{S}$. We assume without loss of generality that the state space $\mathcal{S}$ of $\mathbf{S}$ is the set $\{1, 2, \ldots, M\}$. Let $A = [a_{ij}]_{i,j=1,\ldots,M}$ denote the $M \times M$ transition matrix of $\mathbf{S}$, where $a_{ij} \triangleq \mathbb{P}(S_{t+1} = j|S_t = i)$. If $A$ is ergodic (i.e., irreducible and aperiodic), then there exists a unique probability distribution $\pi$ on $\mathcal{S}$ such that $\pi = \pi A$ (the *stationary distribution* of $\mathbf{S}$), see, e.g., Section 8 of Billingsley [42]. Because in this paper we deal with two-sided random processes, we assume that $\mathbf{S}$ has been initialized with its stationary distribution at some time sufficiently far away in the past, and can therefore be thought of as a two-sided stationary process. Now consider a discrete-time memoryless channel with input alphabet $\mathcal{S}$ and output (observation) alphabet $\mathcal{X} = \mathbb{R}^d$ for some $d < \infty$. It is specified by a set $\{p(\cdot|s) : s = 1, 2, \ldots, M\}$ of transition densities (with respect to $\mu$, the Lebesgue measure on $\mathbb{R}^d$). The channel output sequence $\mathbf{X}$ is the source of interest.

Let us take as the parameter space $\Lambda \subset \mathbb{R}^{M \times M}$ the set of all $M \times M$ transition matrices $[a_{ij}]$, such that all $a_{ij} > a_0$ for some fixed $a_0 > 0$. For each $\theta = [a_{ij}] \in \Lambda$ and each $n \in \mathbb{N}$, the density $dP_\theta^n/d\mu^n$ is given by

$$p_\theta(x^n) = \sum_{s^n \in \mathcal{S}^n} \prod_{i=1}^n a_{s_{i-1}s_i} p(x_i|s_i),$$

where $a_{s_0 s} \equiv \pi_s$ for every $s \in \mathcal{S}$. We assume that the channel transition densities $p(\cdot|s), s \in \mathcal{S}$, are fixed *a priori*, and do not include them in the parametric description of the sources. We do require, though, that

$$\sum_{s \in \mathcal{S}} p(x|s) > 0, \qquad \forall x \in \mathcal{X}$$

and

$$\mathbb{E}_\theta\left\{\log \sum_{s \in \mathcal{S}} p(X|s)\right\} < \infty, \qquad \forall \theta \in \Lambda.$$

We now proceed to verify that Conditions 1–3 of Section III are met.

Let $p_{ij}^{(n)} = \mathbb{P}(S_{t+n} = j|S_t = i)$ denote the $n$-step transition probability for states $i, j \in \mathcal{S}$. The positivity of $A$ implies that the Markov chain $\mathbf{S}$ is *geometrically ergodic*, i.e.,

$$|p_{ij}^{(n)} - \pi_j| \leq C\gamma^n, \qquad \forall i, j \in \mathcal{S}; \forall n \in \mathbb{N} \qquad (27)$$

where $C \geq 0$ and $0 \leq \gamma < 1$, see Theorem 8.9 of Billingsley [42]. Note that (27) implies that

$$d(p^{(n)}(\cdot|i), \pi) \leq MC\gamma^n.$$

This in turn implies that the sequence $\mathbf{S} = \{S_i\}$ is exponentially $\beta$-mixing, see Theorem 3.10 of Vidyasagar [24]. Now, one can show (see Section 3.5.3 of Vidyasagar [24]) that there exists a measurable mapping $F : \mathcal{S} \times [0, 1] \to \mathcal{X}$, such that $X_i = F(S_i, U_i)$, where $\mathbf{U} = \{U_i\}$ is an i.i.d. sequence of random variables distributed uniformly on $[0, 1]$, independently of $\mathbf{S}$. It is not hard to show that, if $\mathbf{S}$ is exponentially $\beta$-mixing, then so is the bivariate process $\{(S_i, U_i)\}$. Finally, because $X_i$ is given by a time-invariant deterministic function of $(S_i, U_i)$, the $\beta$-mixing coefficients of $\mathbf{X}$ are bounded by the corresponding $\beta$-mixing coefficients of $(\mathbf{X}, \mathbf{U})$, and so $\mathbf{X}$ is exponentially $\beta$-mixing as well. Thus, for each $\theta \in \Lambda$, there exists a $\gamma = \gamma(\theta) \in [0, 1)$, such that $\beta_\theta(k) = O(\gamma^k)$, and consequently Condition 1 holds.

To show that Condition 2 holds, we again examine the asymptotic behavior of the Fisher information matrix $I_n(\theta)$ as $n \to \infty$. Under our assumptions on the state transition matrices in $\Lambda$ and on the channel transition densities $\{p(\cdot|s) : s \in \mathcal{S}\}$, we can invoke the results of Section 6.2 in Douc, Moulines and Rydén [43] to conclude that the asymptotic Fisher information matrix $I(\theta) = \lim_{n\to\infty} I_n(\theta)$ exists (though it is not necessarily nonsingular). Thus, Condition 2 is satisfied.

Finally we check Condition 3. The Yatracos class $\mathcal{A}_n$

consists of all sets of the form

$$A_{\theta,\theta'} = \left\{ x^n \in \mathcal{X}^n : \sum_{s^n \in \mathcal{S}^n} \left( \prod_{i=1}^n a_{s_{i-1}s_i} - \prod_{i=1}^n a'_{s_{i-1}s_i} \right) \right.$$
$$\left. \times \prod_{j=1}^n p(x_j|s_j) > 0 \right\}$$

for all $\theta = [a_{ij}], \theta' = [a'_{ij}] \in \Lambda$. The condition $x^n \in A_{\theta,\theta'}$ can be written as $\Pi(x^n, \theta, \theta') > 0$, where for each $x^n$, $\Pi(x^n, \theta, \theta')$ is a polynomial of degree $n$ in the $2M^2$ parameters $a_{ij}, a'_{kl}$, $1 \le i, j, k, l \le M$. Thus, Lemma 2.2 implies that $\mathsf{V}(\mathcal{A}_n) \le 4M^2 \log(4en)$, so Condition 3 is satisfied as well.

## VI. CONCLUSIONS AND FUTURE DIRECTIONS

We have shown that, given a parametric family of stationary mixing sources satisfying some regularity conditions, there exists a universal scheme for joint lossy compression and source identification, with the $n$th-order Lagrangian redundancy and the variational distance between $n$-dimensional marginals of the true and the estimated source both converging to zero as $\sqrt{V_n \log n/n}$, as the block length $n$ tends to infinity. The sequence $\{V_n\}$ quantifies the learnability of the $n$-dimensional marginals. This generalizes our previous results from [9], [10] for i.i.d. sources.

We can outline some directions for future research.

- Both in our earlier work [9], [10] and in the present paper, we assume that the dimension of the parameter space is known *a priori*. It would be of interest to consider the case when the parameter space is finite-dimensional, but its dimension is not known. Thus, we would have a hierarchical model class $\bigcup_{k=1}^\infty \{P_\theta : \theta \in \Lambda^{(k)}\}$, where, for each $k$, $\Lambda^{(k)}$ is an open subset of $\mathbb{R}^k$, and we could use a complexity regularization technique, such as "structural risk minimization" (see, e.g., Lugosi and Zeger [44] or Chapter 6 of Vapnik [22]), to adaptively trade off the estimation and the approximation errors.
- The minimum-distance density estimator of Devroye and Lugosi [34], [35], which plays the key role in our scheme both here and in [9], [10], is not easy to implement in practice, especially for multidimensional alphabets. On the other hand, there are two-stage universal schemes, such as that of Chou, Effros and Gray [14], which do not require memory and select the second-stage code based on pointwise, rather than average, behavior of the source. These schemes, however, are geared toward compression, and do not emphasize identification. It would be worthwhile to devise practically implementable universal schemes that strike a reasonable compromise between these two objectives.
- Finally, neither here nor in our earlier work [9], [10] have we considered the issues of optimality. It would be of interest to obtain lower bounds on the performance of any universal scheme for joint lossy compression and identification, say, in the spirit of minimax lower bounds in statistical learning theory (cf., e.g., Chapter 14 of Devroye, Györfi and Lugosi [21]).

Conceptually, our results indicate that links between statistical modeling (parameter estimation) and universal source coding, exploited in the lossless case by Rissanen [2], [3], are present in the domain of lossy coding as well. We should also mention that another modeling-based approach to universal lossy source coding, due to Kontoyiannis and others (see, e.g., Madiman and Kontoyiannis [45] and references therein), treats code selection as a statistical estimation problem over a class of model distributions in the *reproduction space*. This approach, while closer in spirit to Rissanen's Minimum Description Length (MDL) principle [46], does not address the problem of joint source coding and identification, but it provides a complementary perspective on the connections between lossy source coding and statistical modeling.

## APPENDIX

In this Appendix, we detail some properties of Lagrange-optimal variable-rate vector quantizers. Our exposition is patterned on the work of Linder [19], with appropriate modifications.

As elsewhere in the paper, let $\mathcal{X}$ be the source alphabet and $\widehat{\mathcal{X}}$ the reproduction alphabet, both assumed to be Polish spaces. As before, let the distortion function $\rho$ be induced by a $\rho_{\max}$-bounded metric on a Polish metric space $\mathcal{Y}$ containing $\mathcal{X} \cup \widehat{\mathcal{X}}$. For every $n = 1, 2, \ldots$, define the metric $\rho_n$ on $\mathcal{Y}^n$ by

$$\rho_n(y^n, u^n) \triangleq \frac{1}{n} \sum_{i=1}^n \rho(y_i, u_i).$$

For any pair $P^{(1)}, P^{(2)}$ of probability measures on $\mathcal{X}^n$, let $\mathcal{P}_n(P^{(1)}, P^{(2)})$ be the set of all probability measures on $\mathcal{X}^n \times \widehat{\mathcal{X}}^n$ having $P^{(1)}$ and $P^{(2)}$ as marginals, and define the *Wasserstein metric*

$$\overline{\rho}_n(P^{(1)}, P^{(2)}) \triangleq \inf_{P \in \mathcal{P}_n(P^{(1)}, P^{(2)})} \mathbb{E}_P \{\rho_n(X^n, Y^n)\}$$
$$\equiv \inf_{P \in \mathcal{P}_n(P^{(1)}, P^{(2)})} \int \rho_n(x^n, y^n) dP(x^n, y^n)$$

(See Gray, Neuhoff and Shields [47] for more details and applications.) Note that, because $\rho$ is a bounded metric,

$$\int \rho_n(x^n, y^n) dP(x^n, y^n) \le \rho_{\max} \int 1_{\{x^n \ne y^n\}} dP(x^n, y^n)$$

for all $P \in \mathcal{P}_n(P^{(1)}, P^{(2)})$. Taking the infimum of both sides over all $P \in \mathcal{P}_n(P^{(1)}, P^{(2)})$ and observing that

$$d(P^{(1)}, P^{(2)}) = 2 \inf_{P \in \mathcal{P}_n(P^{(1)}, P^{(2)})} \int 1_{\{x^n \ne y^n\}} dP(x^n, y^n),$$

see, e.g., Section I.5 of Lindvall [48], we get the useful bound

$$\overline{\rho}_n(P^{(1)}, P^{(2)}) \le \frac{1}{2} \rho_{\max} d(P^{(1)}, P^{(2)}). \quad (A.1)$$

Now, for each $n$, let $\mathcal{M}_n$ denote the set of all discrete probability distributions on $\widehat{\mathcal{X}}^n$ with finite entropy. That is, $Q \in \mathcal{M}_n$ if and only if it is concentrated on a finite or a countable set $\{y_i\}_{i \in \mathcal{I}_Q} \subset \widehat{\mathcal{X}}^n$, and

$$H(Q) \triangleq -\sum_{i \in \mathcal{I}_Q} Q(y_i) \log Q(y_i) < \infty.$$

For every $Q \in \mathcal{M}_n$, consider the set $\mathcal{C}(Q)$ of all one-to-one maps $c : \mathcal{I}_Q \to \{0,1\}^*$, such that, for each $c \in \mathcal{C}(Q)$, the collection $\{c(i)\}_{i \in \mathcal{I}_Q}$ satisfies the Kraft inequality, and let

$$\ell_Q \triangleq \min_{c \in \mathcal{C}(Q)} \sum_{i \in \mathcal{I}_Q} \ell(c(i)) Q(y_i)$$

be the minimum expected code length. Since the entropy of $Q$ is finite, there is always a minimizing $c_Q^*$, and the Shannon–Fano bound (see Section 5.4 of Cover and Thomas [1]) guarantees that $\ell_Q \leq H(Q) + 1 < \infty$.

Now, for any $\lambda > 0$, any probability distribution $P$ on $\mathcal{X}^n$, and any $Q \in \mathcal{M}_n$, define

$$L_n(P, Q; \lambda) \triangleq \overline{\rho}_n(P, Q) + \lambda n^{-1} \ell_Q.$$

To give an intuitive meaning to $L_n(P, Q; \lambda)$, let $X^n$ and $Y$ be jointly distributed random variables with $X^n \sim P$ and $Y \sim Q$, such that their joint distribution $\overline{P} \in \mathcal{P}_n(P, Q)$ achieves $\overline{\rho}_n(P, Q)$. Then $L_n(P, Q; \lambda)$ is the expected Lagrangian performance, at Lagrange multiplier $\lambda$, of a *stochastic* variable-rate quantizer which encodes each point $x^n \in \mathcal{X}^n$ as a binary codeword with length $c_Q^*(i)$ and decodes it to $y_i$ in the support of $Q$ with probability $\overline{P}(Y = y_i | X^n = x^n)$.

The following lemma shows that deterministic quantizers are as good as random ones:

*Lemma A.1:* Let $L_P(C^n, \lambda)$ be the expected Lagrangian performance of an $n$-block variable rate quantizer operating on $X^n \sim P$, and let $\widehat{L}_P^n(\lambda)$ be the expected Lagrangian performance, with respect to $P$, of the best $n$-block variable-rate quantizer. Then

$$\widehat{L}_P^n(\lambda) = \inf_{Q \in \mathcal{M}_n} L_n(P, Q; \lambda).$$

*Proof:* Consider any quantizer $C^n = (f, \varphi)$ with $L_P(C^n, \lambda) < \infty$. Let $Q_{C^n}$ be the distribution of $C^n(X^n)$. Clearly, $Q_{C^n} \in \mathcal{M}_n$, and

$$\begin{aligned}
L_P(C^n, \lambda) &= \mathbb{E}\{\rho_n(X^n, C^n(X^n))\} + \lambda\, \mathbb{E}\{\ell_n(f(X^n))\} \\
&\geq \overline{\rho}_n(P, Q_{C^n}) + \lambda n^{-1} \ell_{Q_{C^n}} \\
&= L_n(P, Q_{C^n}; \lambda).
\end{aligned}$$

Hence, $\widehat{L}_P^n(\lambda) \geq \inf_{Q \in \mathcal{M}_n} L_n(P, Q; \lambda)$. To prove the reverse inequality, suppose that $X^n \sim P$ and $Y \sim Q$ achieve $\overline{\rho}_n(P, Q)$ for some $Q \in \mathcal{M}_n$. Let $\overline{P}$ be their joint distribution. Let $\{y_i\}_{i \in \mathcal{I}_Q} \subset \widehat{\mathcal{X}}^n$ be the support of $Q$, let $c_Q^* : \mathcal{I}_Q \to \{0,1\}^*$ achieve $\ell_Q$, and let $\mathcal{S} = \{c_Q^*(i)\}_{i \in \mathcal{I}_Q}$ be the associated binary code. Define the quantizer $\widetilde{C}^n = (f, \varphi)$ by

$$\varphi(s) = y_i \quad \text{if } s = c_Q^*(i)$$

and

$$f(x^n) = \arg\min_{s \in \mathcal{S}} \left(\rho_n(x^n, \varphi(s)) + \lambda \ell_n(s)\right).$$

Then

$$L_P(C^n, \lambda) = \mathbb{E}_P\left\{\min_{s \in \mathcal{S}}\left(\rho_n(X^n, \varphi(s)) + \lambda \ell_n(s)\right)\right\}.$$

On the other hand,

$$\begin{aligned}
&L_n(P, Q; \lambda) \\
&= \mathbb{E}_{\overline{P}}\left\{\overline{\rho}_n(X_1^n, Y) + \lambda n^{-1} \ell_Q\right\} \\
&= \int dP(x^n) \sum_{i \in \mathcal{I}_Q} \left(\rho_n(x^n, y_i) + \lambda \ell_n(c_Q^*(i))\right) \\
&\qquad \times \overline{P}(Y = y_i | X^n = x^n) \\
&\geq \int dP(x^n) \min_{i \in \mathcal{I}_Q} \left(\rho_n(x^n, y_i) + \lambda \ell_n(c_Q^*(i))\right) \\
&= \int dP(x^n) \min_{s \in \mathcal{S}} \left(\rho_n(x^n, \varphi(s)) + \lambda \ell_n(s)\right) \\
&\equiv L_P(C^n, \lambda),
\end{aligned}$$

so that $\inf_{Q \in \mathcal{M}_n} L_n(P, Q; \lambda) \geq \widehat{L}_P^n(\lambda)$, and the lemma is proved. $\blacksquare$

The following lemma gives a useful upper bound on the Lagrangian mismatch:

*Lemma A.2:* Let $P, P'$ be probability distributions on $\mathcal{X}^n$. Then

$$\left|\widehat{L}_P^n(\lambda) - \widehat{L}_{P'}^n(\lambda)\right| \leq \frac{1}{2}\rho_{\max} d(P, P').$$

*Proof:* Suppose $\widehat{L}_P^n(\lambda) \geq \widehat{L}_{P'}^n(\lambda)$. Let $Q'$ achieve $\inf_{Q \in \mathcal{M}_n} L_n(P', Q; \lambda)$ (or be arbitrarily close). Then

$$\begin{aligned}
&\widehat{L}_P^n(\lambda) - \widehat{L}_{P'}^n(\lambda) \\
&\overset{(a)}{=} \inf_{Q \in \mathcal{M}_n} L_n(P, Q; \lambda) - \inf_{Q \in \mathcal{M}_n} L_n(P', Q; \lambda) \\
&= \inf_{Q \in \mathcal{M}_n} L_n(P, Q; \lambda) - L_n(P', Q'; \lambda) \\
&\leq L_n(P, Q'; \lambda) - L_n(P', Q'; \lambda) \\
&\overset{(b)}{=} \overline{\rho}_n(P, Q') + \lambda n^{-1} \ell_{Q'} - \overline{\rho}_n(P', Q') - \lambda n^{-1} \ell_{Q'} \\
&= \overline{\rho}_n(P, Q') - \overline{\rho}_n(P', Q') \\
&\overset{(c)}{\leq} \overline{\rho}_n(P, P') \\
&\overset{(d)}{\leq} \frac{1}{2}\rho_{\max} d(P, P'),
\end{aligned}$$

where in (a) we used Lemma A.1) in (b) we used the definition of $L_n(\cdot, Q'; \lambda)$, in (c) we used the fact that $\overline{\rho}_n$ is a metric and the triangle inequality, and in (d) we used the bound (A.1). $\blacksquare$

Finally, the lemma below shows that, for bounded distortion functions, Lagrange-optimal quantizers have finite codebooks:

*Lemma A.3:* For positive integers $N, L$, let $\mathcal{Q}_n(N, L)$ denote the set of all zero-memory variable-rate quantizers with block length $n$, such that for every $C^n \in \mathcal{Q}_n(N, L)$, the associated binary code $\mathcal{S}$ of $C^n$ satisfies $|\mathcal{S}| \leq N$ and $\ell(s) \leq L$ for every $s \in \mathcal{S}$. Let $P$ be a probability distribution on $\mathcal{X}^n$. Then

$$\widehat{L}_P^n(\lambda) = \inf_{C^n \in \mathcal{Q}_n(N, L)} L_P(C^n, \lambda),$$

with $N \leq 2^{2n\rho_{\max}/\lambda}$ and $L \leq 2n\rho_{\max}/\lambda$.

*Proof:* Let $C_*^n$ with encoder $f_* : \mathcal{X}^n \to \mathcal{S}$ and decoder $\varphi_* : \mathcal{S} \to \widehat{\mathcal{X}}^n$ achieve the $n$th-order optimum $\widehat{L}_P^n(\lambda)$ for $P$. Let $s_0 \in \mathcal{S}$ be the shortest binary string in $\mathcal{S}$, i.e.,

$$\ell(s_0) = \min_{s \in \mathcal{S}} \ell(s).$$

Without loss of generality, we can take $f_*$ as the minimum-

distortion encoder, i.e.,

$$f_*(x^n) = \arg\min_{s \in \mathcal{S}} \left( \rho_n(x^n, \varphi_*(s)) + \lambda \ell_n(s) \right).$$

Thus, for any $s \in \mathcal{S}$ and any $x^n \in f_*^{-1}(s)$,

$$\rho_n(x^n, \varphi_*(s)) + \lambda \ell_n(s) \leq \rho_n(x^n, \varphi_*(s_0)) + \lambda \ell_n(s_0).$$

Hence, $\ell(s) \leq n\rho_{\max}/\lambda + \ell(s_0)$ for all $s \in \mathcal{S}$. Furthermore, $L_P(C_*^n, \lambda) \geq \lambda \mathbb{E}_P \{\ell_n(f_*(X^n))\} \geq \lambda \ell_n(s_0)$.

Now pick an arbitrary reproduction string $\widehat{x}_0^n \in \widehat{\mathcal{X}}^n$, let $\varepsilon$ be the empty binary string (of length zero), and let $C_0^n$ be the zero-rate quantizer with the constant encoder $f_0(x^n) = \varepsilon$ and the decoder $\varphi_0(\varepsilon) = \widehat{x}_0^n$. Then $L_P(C_0^n, \lambda) = \mathbb{E}_P\{\rho_n(X^n, \widehat{x}_0^n)\} + \lambda \ell_n(\varepsilon) \leq \rho_{\max}$. On the other hand, $L_P(C_*^n, \lambda) \leq L_P(C_0^n, \lambda)$. Therefore,

$$\lambda \ell_n(s_0) \leq L_P(C_*^n, \lambda) \leq L_P(C_0^n, \lambda) \leq \rho_{\max},$$

so that $\ell(s_0) \leq n\rho_{\max}/\lambda$. Hence,

$$\ell(s) \leq 2n\rho_{\max}/\lambda, \qquad \forall s \in \mathcal{S},$$

Since the strings in $\mathcal{S}$ must satisfy Kraft's inequality, we have

$$1 \geq \sum_{s \in \mathcal{S}} 2^{-\ell(s)} \geq |\mathcal{S}| 2^{-2n\rho_{\max}/\lambda},$$

which implies that $|\mathcal{S}| \leq 2^{2n\rho_{\max}/\lambda}$. ∎

## ACKNOWLEDGMENT

## REFERENCES

[1] T. M. Cover and J. A. Thomas, *Elements of Information Theory*. New York: Wiley, 1991.

[2] J. Rissanen, "Universal coding, information, prediction, and estimation," *IEEE Trans. Inform. Theory*, vol. IT-30, no. 4, pp. 629–636, July 1984.

[3] ——, "Fisher information and stochastic complexity," *IEEE Trans. Inform. Theory*, vol. 42, no. 1, pp. 40–47, January 1996.

[4] J. Ziv and A. Lempel, "Compression of individual sequences by variable-rate coding," *IEEE Trans. Inform. Theory*, vol. IT-24, pp. 530–536, September 1978.

[5] J. C. Kieffer, "Strongly consistent code-based identification and order estimation for constrained finite-state model classes," *IEEE Trans. Inform. Theory*, vol. 39, no. 3, pp. 893–902, May 1993.

[6] N. Merhav, "Bounds on achievable convergence rates of parameter estimation via universal coding," *IEEE Trans. Inform. Theory*, vol. 40, no. 4, pp. 1210–1215, July 1994.

[7] T. Weissman and E. Ordentlich, "The empirical distribution of rate-constrained source codes," *IEEE Trans. Inform. Theory*, vol. 51, no. 11, pp. 3718–3733, November 2005.

[8] G. Tao, *Adaptive Control Design and Analysis*. Hoboken: Wiley, 2003.

[9] M. Raginsky, "Joint fixed-rate universal lossy coding and identification of continuous-alphabet memoryless sources," *IEEE Trans. Inform. Theory*, vol. 54, no. 7, pp. 3059–3077, July 2008.

[10] ——, "Joint universal lossy coding and identification of i.i.d. vector sources," in *Proc. IEEE Int. Symp. on Information Theory*, Seattle, July 2006, pp. 577–581.

[11] E.-H. Yang and Z. Zhang, "On the redundancy of lossy source coding with abstract alphabets," *IEEE Trans. Inform. Theory*, vol. 45, no. 4, pp. 1092–1110, May 1999.

[12] V. N. Vapnik and A. Y. Chervonenkis, "On the uniform convergence of relative frequencies of events to their probabilities," *Theory Probab. Appl.*, vol. 16, pp. 264–280, 1971.

[13] P. A. Chou, T. Lookabaugh, and R. M. Gray, "Entropy-constrained vector quantization," *IEEE Trans. Acoust. Speech Signal Processing*, vol. 37, no. 1, pp. 31–42, January 1989.

[14] P. A. Chou, M. Effros, and R. M. Gray, "A vector quantization approach to universal noiseless coding and quantization," *IEEE Trans. Inform. Theory*, vol. 42, no. 4, pp. 1109–1138, July 1996.

[15] R. M. Gray, *Entropy and Information Theory*. New York: Springer-Verlag, 1990.

[16] V. A. Volkonskii and Y. A. Rozanov, "Some limit theorems for random functions, I," *Theory Probab. Appl.*, vol. 4, pp. 178–197, 1959.

[17] ——, "Some limit theorems for random functions, II," *Theory Probab. Appl.*, vol. 6, pp. 186–198, 1961.

[18] D. L. Neuhoff and R. K. Gilbert, "Causal source codes," *IEEE Trans. Inform. Theory*, vol. IT-28, no. 5, pp. 701–713, September 1982.

[19] T. Linder, "Learning-theoretic methods in vector quantization," in *Principles of Nonparametric Learning*, L. Györfi, Ed. New York: Springer-Verlag, 2001.

[20] M. Effros, P. A. Chou, and R. M. Gray, "Variable-rate source coding theorems for stationary nonergodic sources," *IEEE Trans. Inform. Theory*, vol. 40, no. 6, pp. 1920–1925, November 1994.

[21] L. Devroye, L. Györfi, and G. Lugosi, *A Probabilistic Theory of Pattern Recognition*. New York: Springer-Verlag, 1996.

[22] V. N. Vapnik, *Statistical Learning Theory*. New York: Wiley, 1998.

[23] L. Devroye and G. Lugosi, *Combinatorial Methods in Density Estimation*. New York: Springer-Verlag, 2001.

[24] M. Vidyasagar, *Learning and Generalization*, 2nd ed. London: Springer-Verlag, 2003.

[25] N. Sauer, "On the density of families of sets," *J. Combin. Theory Ser. A*, vol. 13, pp. 145–147, 1972.

[26] R. M. Dudley, "Central limit theorems for empirical measures," *Ann. Probab.*, vol. 6, pp. 898–929, 1978.

[27] M. Karpinski and A. Macintyre, "Polynomial bounds for VC dimension of sigmoidal and general Pfaffian neural networks," *J. Comput. Sys. Sci.*, vol. 54, pp. 169–176, 1997.

[28] S. N. Bernstein, "Sur l'extension du théorème limite du calcul des probabilités aux sommes de quantités dependantes," *Math. Ann.*, vol. 97, pp. 1–59, 1927.

[29] B. Yu, "Rates of convergence for empirical processes of stationary mixing sequences," *Ann. Probab.*, vol. 22, no. 1, pp. 94–116, 1994.

[30] R. Meir, "Nonparametric time series prediction through adaptive model selection," *Machine Learning*, vol. 39, pp. 5–34, 2000.

[31] A. Mokkadem, "Mixing properties of ARMA processes," *Stochastic Process. Appl.*, vol. 29, pp. 309–315, 1988.

[32] B. S. Clarke and A. R. Barron, "Information-theoretic asymptotics of Bayes methods," *IEEE Trans. Inform. Theory*, vol. 36, no. 3, pp. 453–471, May 1990.

[33] Y. G. Yatracos, "Rates of convergence of minimum distance estimates and Kolmogorov's entropy," *Ann. Math. Statist.*, vol. 13, pp. 768–774, 1985.

[34] L. Devroye and G. Lugosi, "A universally acceptable smoothing factor for kernel density estimation," *Ann. Statist.*, vol. 24, pp. 2499–2512, 1996.

[35] ——, "Nonasymptotic universal smoothing factors, kernel complexity and Yatracos classes," *Ann. Statist.*, vol. 25, pp. 2626–2637, 1997.

[36] P. Elias, "Universal codeword sets and representations of the integers," *IEEE Trans. Inform. Theory*, vol. IT-21, no. 2, pp. 194–203, March 1975.

[37] L. Devroye and L. Györfi, "Distribution and density estimation," in *Principles of Nonparametric Learning*, L. Györfi, Ed. New York: Springer-Verlag, 2001.

[38] I. Kontoyiannis and J. Zhang, "Arbitrary source models and Bayesian codebooks in rate-distortion theory," *IEEE Trans. Inform. Theory*, vol. 48, no. 8, pp. 2276–2290, August 2002.

[39] A. Klein and P. Spreij, "The Bezoutian, state space realizations and Fisher's information matrix of an ARMA process," *Lin. Algebra Appl.*, vol. 416, pp. 160–174, 2006.

[40] P. J. Bickel, Y. Ritov, and T. Rydén, "Asymptotic normality of the maximum-likelihood estimator for general hidden Markov models," *Ann. Statist.*, vol. 26, no. 4, pp. 1614–1635, 1997.

[41] Y. Ephraim and N. Merhav, "Hidden Markov processes," *IEEE Trans. Inform. Theory*, vol. 48, no. 6, pp. 1518–1569, June 2002.

[42] P. Billingsley, *Probability and Measure*, 3rd ed. Wiley, New York.

[43] R. Douc, É. Moulines, and T. Rydén, "Asymptotic properties of the maximum likelihood estimator in autoregressive models with Markov regime," *Ann. Statist.*, vol. 32, no. 5, pp. 2254–2304, 2004.

[44] G. Lugosi and K. Zeger, "Concept learning using complexity regularization," *IEEE Trans. Inform. Theory*, vol. 42, no. 1, pp. 48–54, January 1996.

[45] M. Madiman and I. Kontoyiannis, "Second-order properties of lossy likelihoods and the MLE/MDL dichotomy in lossy compression," Brown University, APPTS Report No. 04-5, May 2004, available [Online] at http://www.dam.brown.edu/ptg/REPORTS/04-5.pdf.

[46] A. Barron, J. Rissanen, and B. Yu, "Minimum description length principle in coding and modeling," *IEEE Trans. Inform. Theory*, vol. 44, no. 6, pp. 2743–2760, October 1998.

[47] R. M. Gray, D. L. Neuhoff, and P. S. Shields, "A generalization of Ornstein's $\bar{d}$ distance with applications to information theory," *Ann. Probab.*, vol. 3, no. 2, pp. 315–328, 1975.

[48] T. Lindvall, *Lectures on the Coupling Method*. New York: Dover, 2002.